

A Survey on Representing Linguistic Style: Challenges and Opportunities

Anonymous ACL submission

Abstract

Although representation learning has transformed semantic modeling in NLP, representations of linguistic style remain underexplored—partly due to conflicting definitions of style within and beyond NLP or unclear immediate advantages of separate style representations. In this survey, we provide an overview of style conceptualizations across different research fields with a focus on NLP and (socio-)linguistics and suggest a working definition of style for practitioners. Then, we review methods for creating and evaluating style representations. We conclude by discussing how style representations can make crucial contributions to the modern NLP pipeline (e.g., in dataset curation or evaluation) and to the application of NLP methods in other fields. Throughout our survey, we sketch pressing open research questions in the landscape of style representations, emphasizing the need for better evaluation approaches and more comprehensive style representations.

1 Introduction

The Lego Grad Student¹ posted in July 2020, *Videoconferencing from his apartment with his advisor, the grad student feels like the victim of a home invasion.*

Now consider a rephrasing by GPT-5.2 using the Wikipedia-style prompt from Maini et al. (2024):

While conducting a videoconference with his academic advisor from his apartment, the graduate student experiences the interaction as an intrusion into his private living space.

The linguistic style of the original post (e.g., more informal, compact) likely contributed to the 3k likes it received. Style can affect a reader’s perception as

¹The “Lego Grad Student” is an online creator that received engagement on Twitter and Instagram with photos of LEGO figures playing out scenes in a grad student’s life. This message was posted during the COVID-19 pandemic.

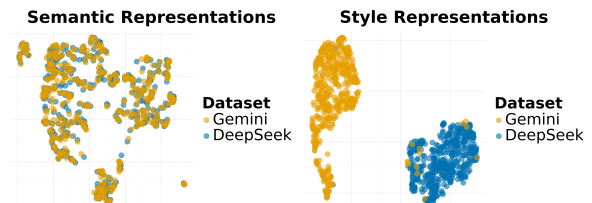


Figure 1: **Semantic representations differ from style representations** We compare reasoning traces from Muennighoff et al. (2025), generated with Gemini and DeepSeek for the same reasoning problems. Style representations can distinguish between the two models—confirming results in Rivera Soto et al. (2023)—while semantic representations overlap. See §B.1 for details.

it can, for example, influence engagement (Munaro et al., 2024; Banerjee and Urminsky, 2025) and change the persuasiveness of arguments (El Baff et al., 2020; Parhankangas and Renko, 2017). Moreover, style also influences human quality ratings of generated content² (Cai et al., 2024; Wu and Aji, 2025) leading to one of our main takeaways: *If you care about LLMs, then style matters.*

However, style is often disregarded in NLP. As a result, language models can be brittle across (or not robust to) style-like features: rephrasing prompts in different styles leads to different performances (Mizrahi et al., 2024; Wahle et al., 2024); LLM judges can prefer long, formal, or synthetic texts over relevance (Cao, 2025; Feuer et al., 2025; Wu and Aji, 2025); machine translations sound older and more male than the original (Hovy et al., 2020); models are biased against non-majority varieties (Fleisig et al., 2024; Hofmann et al., 2024; Liang et al., 2023) and perform worse on non-standard spellings (Ebrahimi et al., 2018; Li et al., 2019), simple and informal styles, and registers like poetry (Anschütz et al., 2025; Cao, 2025; Qi et al., 2021; Zhao et al., 2025). This brittleness might increase as we train on more synthetic data (Guo et al., 2024).





²People might prefer the style of certain LLMs over others, e.g., Claude’s style over ChatGPT’s, or want to customize an LLM’s style (OpenAI, 2025); see Personalization in §5.2.

063 *Style representations* (i.e., vectors with entries
064 that are optimized for style information) can help:
065 They can improve model robustness by supporting
066 the curation of stylistically diverse (post-)training
067 datasets, support text generation in and evaluate
068 adherence to a target style, help machine text de-
069 tection, and enable new tasks (e.g., retrieval of doc-
070 uments in a target style). In the social sciences and
071 humanities, they can support the analysis of literary
072 texts and style dynamics in dialogue. A detailed
073 discussion of these possibilities follows in §5.

074 Semantic representations are often also sensitive
075 to style information as word prediction tasks also
076 need style information (Nguyen and Grieve, 2020;
077 Goldberg, 2019; Tenney et al., 2018; Miaschi et al.,
078 2020; Wegmann and Nguyen, 2021). However, we
079 believe that semantic representations alone are in-
080 sufficient for modeling style: They are usually not
081 evaluated on style-related tasks (Enevoldsen et al.,
082 2025; Muennighoff et al., 2023) and have limited
083 sensitivity to style (e.g., Mickus and Copot, 2024;
084 Zhang et al., 2023b). Most importantly, they are
085 trained to focus primarily on semantic information,
086 making it difficult to investigate the style of texts
087 separately from content (cf. Fig. 1).

088 In contrast to semantic representations, only a
089 few community-vetted and broadly tested methods
090 exist for representing the style of texts. **The main
091 goals of this paper are** to promote the wider adop-
092 tion of linguistic style representations within and
093 beyond NLP, guide practitioners towards key re-
094 sources, and highlight key challenges and research
095 directions in the study of style representations.

096 **With this paper, we contribute:**

- 097 • an overview of style definitions in linguistics
098 and NLP, including our own definition (§2)
- 099 • an overview of methods for representing (§3)
100 and evaluating (§4) style representations
- 101 • a discussion of why style representations are
102 useful for modern NLP and other fields (§5)
- 103 •  practical resources,  open research ques-
104 tions, and  calls to action (several sections)
- 105 •  an expanding GitHub repository³ collect-
106 ing datasets, tools, and other resources

107 Despite the significant attention given to style in
108 other modalities (e.g., speech), text-based NLP has
109 lagged behind, highlighting the need for this survey.
110 In line with this limited coverage, most of the work
111 we discuss focuses on English texts, but we urge

³<https://anonymous.4open.science/w/StyleSurvey-60B0/>

the NLP community to consider more languages
and modalities in the future.

2 Style conceptualizations

Linguists often define style as a distinctive
pattern in language for some object of study
(e.g., for an author or group), while NLP
researchers often use “style” more loosely.

2.1 Style in linguistics

115 Researchers working with style often aim to de-
116 scribe a text’s structural linguistic features (i.e.,
117 how something is said) more so than its seman-
118 tic meaning⁴ (i.e., what is said). However, some
119 linguistics researchers might disagree with such a
120 separation (see §C), finding that style and content
121 are intertwined, at least to some extent (cf. §2.3).
122 Studying style might then be understood as study-
123 ing what makes a phrasing distinctive within a set
124 of possibilities (Irvine, 2001), for instance, how
125 speakers use linguistic choices related to external
126 factors like social background, identity or register.
127 Overall, we emphasize that *style is an elusive term
128 that has been defined in many different, sometimes
129 inconsistent, ways in linguistics and other fields.*⁵
130

What are the objects of study? Style is usu-
ally studied in a relative sense, as a distinctive dif-
ference between objects of study (Irvine, 2001);
however, these objects vary. In (socio-)linguistics,
style has often been discussed as inter-individual
variation—the idiosyncratic choices that poten-
tially distinguish individuals from each other,
often referred to as their *idiolect* (Coulthard,
2004)—and intra-individual variation (Bell, 1984;
Irvine, 2001; Labov, 2006; Meyerhoff, 2006; Wag-
ner, 2025)—the change in the same speaker’s lan-
guage across situations. Famously, Labov (1972)
discovered that individuals’ speaking style becomes
more formal as they pay more attention to their
speech and more casual as they pay less attention.
Sociolinguists have additionally studied style as
inter-group variation—differences in the language
of people identifying with different social groups
(Bell, 1984; Eckert, 2008; Irvine, 2001; Kristiansen,
2024). For example, *g*-dropping (*going* vs. *goin*’)

⁴Or: referential meaning (Campbell-Kibler, 2011; Labov,
1972; Lavandera, 1978; Nguyen et al., 2016, 2021). Two
variants have the same referential meaning if they are the
same in a truth-conditional sense (i.e., true in exactly the same
situations), while the “social” or “stylistic significance” might
differ considerably (Labov, 1972; Weiner and Labov, 1983).

⁵See §C for an overview of other areas interested in style.

152 may indicate a person’s association with a southern
153 U.S. region (Campbell-Kibler, 2007).

154 Domains or registers have also been objects of
155 style research (Biber and Conrad, 2019; Grieve,
156 2023). Literature from a historical period, nov-
157 els by a specific author, news reports, blogs, and
158 conversations can display very different linguistic
159 patterns, which might be called the style of that
160 historical period, literary author, news report, blog,
161 or conversation (Biber and Conrad, 2019; Grieve
162 et al., 2011; Hicke and Mimno, 2025; Irvine, 2001).

163 Researchers have considered more objects of
164 study than we discuss, like the communication envi-
165 ronment (e.g., speech before a crowd or a courtroom
166 in Ervin-Tripp, 2001) or the communicative manner
167 (spontaneous vs. read speech in Williams and King,
168 2019). Researchers can also study combinations of
169 these objects (e.g., courtroom speeches by one indi-
170 vidual) or an object only in certain contexts (e.g., a
171 social group discussing a certain topic). For exam-
172 ple, Holliday (2021) finds that biracial Black men
173 displayed fewer African American⁶ intonational
174 features when discussing police narratives.

175 **What is the function of style?** Style might also
176 be defined as patterns in language tied to a spe-
177 cific function. Some scholars argue that style is
178 fundamentally embedded in social meaning, index-
179 ing social background and shaping social identity
180 (Campbell-Kibler et al., 2006; Coupland, 2007;
181 Eckert, 2008, 2012). For example, Labov (1972)
182 found that differences in the pronunciation of /r/ cor-
183 related with social class, and Eckert (1989) found
184 that self-identified “burnouts” at a Detroit school
185 used more non-standard linguistic features (e.g.,
186 *gonna*) than college-bound “jocks” (e.g., *going to*).

187 Labov originally viewed a speaker’s vernacular
188 as a reflection of their social identity, not an active
189 choice (Labov, 1972). More recent sociolinguistic
190 approaches see style as more agentive—not only re-
191 flecting identity but also performing and construct-
192 ing it (Eckert, 2012). For example, the development
193 of linguistic practices of trans activists can be tied
194 to their agency in creating identity (Zimman, 2019),
195 and speakers may choose styles for performative
196 functions like getting attention (Ervin-Tripp, 2001).

⁶While several linguistic features can describe both styles and dialects, dialects are typically not called styles but distinct types of language variation more clearly tied to speakers’ social backgrounds and geographic regions (Biber and Conrad, 2019; Grieve et al., 2025). Nonetheless, some researchers also consider dialects as a kind of social style (Coupland, 2007). We do not specifically exclude dialects in our definition (§2.3), but our focus remains on non-dialectal stylistic variation.

197 Style can serve communicative functions in an
198 interaction (Coupland, 2007): Speakers may align
199 with (accommodate) or distance themselves from
200 the style of interlocutors or audiences (Bell, 1984;
201 Giles and Powesland, 1975; Giles et al., 1991;
202 Khaleghzadegan et al., 2024), thereby shaping so-
203 cial relationships and interactions (Coupland, 2007).
204 For example, Bell (2014) found that New Zealand
205 newscasters shifted their pronunciation when talk-
206 ing to audiences of higher or lower status.

207 Finally, some consider style to be aesthetic, with
208 no or limited function, and instead prefer the term
209 *register* for varieties of language associated with
210 a particular situational context (Biber and Conrad,
211 2019). When considering register as style, style
212 might serve further functions like structuring dis-
213 course and fulfilling communicative purposes.

2.2 Style in NLP 214

215 Some work in NLP uses the term style in ways
216 broadly consistent with linguistics, aiming to study
217 formal/informal styles and literary authorial styles
218 (e.g., Jhamtani et al., 2017; Rao and Tetreault, 2018;
219 Wegmann and Nguyen, 2021); however, others in-
220 creasingly use style as an umbrella term for general
221 attributes of texts that vary across datasets (Jin et al.,
222 2022) such as the sentiment of a text (Reif et al.,
223 2022; Shen et al., 2017), but do not necessarily
224 align with a typical linguistic definition of style.

225 **Separating content and style** As in linguistics
226 (§2.1), work in NLP finds that content and style
227 are often correlated (Jafaritazehjani et al., 2020;
228 Mikros and Argiri, 2007). Still, separating style and
229 content tends to be a natural distinction for many
230 NLP applications. Specifically, NLG systems have
231 to fundamentally determine what information to
232 generate—the knowledge, or message—and what
233 style to generate it in (Gatt and Krahmer, 2018).
234 While neural NLG systems often handle content and
235 style implicitly, generating texts end-to-end without
236 explicit planning stages, the distinction between
237 style and content remains useful in practice, for
238 example, when curating datasets, rephrasing and
239 adapting texts, or evaluating the factual correctness
240 of model outputs (§5).

2.3 A working definition for style in NLP 241

242 We propose a working definition of style for NLP
243 practitioners.⁷ Throughout the paper, we consis-

⁷Our definition does not specifically exclude concepts like dialects, registers, or varieties for practical reasons: (i) the

tently use the same colors for the same concepts.

Definition A linguistic style consists of *distinctive patterns in language use* for an **object of study** (e.g., individuals, a group of authors in a given register) in its **lexical, syntactic, morphological, orthographic, discourse, phonetic, etc. composition**. These patterns should **not chiefly measure**, but can correlate with, **semantic meaning**.

For example, a person discussing American football might talk more casually than when discussing ballroom dance, yet some underlying linguistic features may remain consistent in both situations and carry social meaning (§2), e.g., about the speaker’s upbringing. When studying style, we might study the differences or commonalities between discussing American football and ballroom dance, depending on the object of study, i.e., whether we are currently interested in a specific individual, demographic, situation, etc.

3 Representing style

Linguistic style is usually operationalized with patterns in linguistic features like function words or automatically-learned representations like neural text representations.

3.1 Predefined features

Style is often operationalized as the systematic variation of linguistic features, which can span various linguistic levels including morphology, orthography, syntax, and discourse (Biber and Conrad, 2019; Crystal and Davy, 1969; Grieve, 2007; Kniffka, 2007; Labov, 1972; Neal et al., 2017; Stamatatos, 2009). 🛠️ App. Tab. 1 gives example features (e.g., g-dropping) at each level; 🛠️ §D lists tools for extracting predefined features. The primary appeal of predefined features is that they are supported by linguistic theory, have been tested extensively, and are generally interpretable (i.e., have a meaning understandable to humans). The features can be used with statistical approaches like logistic regression or dimensionality reduction with factor analysis to determine how important each feature is. This transparency is especially important in high-stakes settings, such as forensic linguistics, where separation between such terms is not consistent in linguistics, and (ii) computational style representations are commonly expected to be sensitive to dialect, register, and variety information (§4). We leave further practical disentanglement between style and other terms for future work.

it is crucial to explain a model’s decision-making process (Argamon, 2018; Grant, 2022).

One such feature-based style operationalization is stylometry, which measures the frequencies of linguistic features that help discriminate between author styles. There is no fixed set of features that work for every individual, despite much work attempting to find one (Juola, 2006; Nini, 2023); instead, the features often depend on the nature of the data (e.g., genre, register, amount of data, language) (Argamon, 2018). Nonetheless, function words (i.e., words like prepositions and conjunctions that primarily serve a grammatical role) and character n-grams (i.e., n successive characters), in particular, have proven quite effective at discriminating authors (Grieve, 2007; Houvardas and Stamatatos, 2006; Peng et al., 2003; Kestemont, 2014; Mosteller and Wallace, 1963) and speakers (Aggazzotti and Smith, 2025; Aggazzotti et al., 2024; Doddington, 2001; Sergidou et al., 2023; Tripto et al., 2023).

Other feature operationalizations serve different purposes related to style. For example, Multidimensional Analysis (MDA) (Biber, 1988) is used to determine how texts differ in their communicative function and originally relied on mostly grammatical category-related features (e.g. nouns, verbs); however, modern extensions (e.g., Clarke and Grieve, 2017; Grieve et al., 2011) additionally include more complex features, such as syntactic constructions and semantic classes.

3.2 Automatically-learned features

By automatically-learned features or embeddings, we mean vector representations of text produced by (usually neural) models. In contrast to predefined features, automatically-learned features do not rely on specific, established features but can automatically discover style patterns. Further, they often perform better than predefined features on downstream tasks, but are usually less interpretable. Because it is difficult to operationalize definitions of style, models are usually optimized in proxy downstream tasks, such as authorship verification or style transfer. 🛠️ See §D for links to models.

Authorship verification The most popular approach to date trains models with a contrastive objective (Dong and Shen, 2018; Khosla et al., 2020) to learn representations where two text samples are close together in vector space if they are written by the same author and far apart otherwise (Andrews and Bishop, 2019; Khan et al., 2021; Kim et al.,

2025; Man and Huu Nguyen, 2024; Rivera Soto et al., 2021; Sawatphol et al., 2022; Thakrar et al., 2025; Wang et al., 2023; Wegmann et al., 2022). Representations trained on this task have been shown to capture stylistic information (Wang et al., 2023; Wegmann and Nguyen, 2021).

Since training datasets may contain undesired correlations—for example between style and **content** when an author only writes about one topic—some work creates harder positive (i.e., same author) and negative (i.e., different author) pairs to improve **disentanglement** (Man and Huu Nguyen, 2024; Patel et al., 2025). For example, Wegmann et al. (2022) use negative pairs that are approximately about the same topic, and Patel et al. (2025) leverage LLMs to create a synthetic dataset of near-exact paraphrases by varying predefined features. Building on such disentanglement strategies, recent work generalizes style representations to multilingual settings (Kim et al., 2025; Qiu et al., 2025), where negative pairs must be carefully constructed to avoid trivial cross-lingual differences.

Style transfer Another line of work learns representations via style-transfer, aiming to rewrite text for a stylistic attribute without altering its semantic meaning (Cheng et al., 2020b; John et al., 2019; Shen et al., 2017; Zhu et al., 2024). For instance, a model may be trained to convert formal text into informal text, conditioned on both the input and an embedding of the target style. Under this objective, embeddings learn features indicative of informality.

These methods usually rely on explicit style-content disentanglement and tend to learn representations that are more narrow in scope, often tied to single attributes (e.g., politeness) or differences between two corpora (Shen et al., 2017). John et al. (2019) train an auto-encoder to produce a style and a content vector, imposing a style classification loss on the style representation and an adversarial style classification loss on the content vector. Cheng et al. (2020b) minimizes the estimated mutual information between the style and content representations.

Interpretable LLM-guided stylometry A distinctive method is LISA (Patel et al., 2023), which learns embeddings where each dimension is an *interpretable* feature (e.g., use of an elongated word). The authors create a synthetic dataset by prompting GPT-3 for stylometric features, then train an EncT5 (Liu et al., 2022) model to predict the presence of each feature in a text sample. Because distances in this space are not well-defined, they fit a

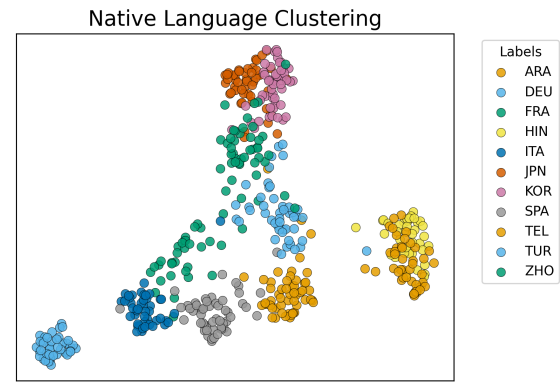


Figure 2: **By-product of authorship verification training** Stylistic representations, though trained on the “idiolectal” authorship verification task, cluster TOEFL (Test of English as a Foreign Language) essays by the native language of the writer. See §B.3 for more details.

linear transformation on the authorship verification task. LISA is the first method to use LLM-based automatic labeling for style representations, offering a middle ground between hand-crafted features derived by human experts and automatically-learned representations. However, limitations of LLMs might need to be considered (§ 1, § 5.1).

3.3 The future of style representations

Define what we want to represent

Training representations on the authorship verification task implicitly defines style as the idiosyncrasies exhibited by authors in certain corpora (Zhu and Jurgens, 2021). However, because the contrastive dataset might never pair authors from the same social group as negatives, a representation may inadvertently primarily encode group-level features. Indeed, we find that various “idiolectal” style representations encode features discriminative of writers’ native language in Fig. 2. We call on the community to explicitly define their object of study (e.g., idiolect, cf. § 2.3), use learning approaches like hard negatives to control for other concepts (e.g., variation within the same social group), and evaluate whether representations primarily capture variations for the defined object of study.

? Build general-purpose style embeddings

It remains an open challenge to learn general-purpose style embeddings that cover as many objects of study as possible and are, for example, sensitive to individual, group, register, and time period variation at the same time. For this, new training objectives could explicitly target different objects of study. Using multiple objectives might

require stronger disentanglement objectives, for example, based on minimizing mutual information of two representations (Cheng et al., 2020a), adversarial objectives (John et al., 2019), or by employing VAEs to explicitly disentangle between syntax and semantics (Chen et al., 2019; Bao et al., 2019).

? Improve training

There are several other areas in training that remain underexplored. For example, fine-tuning newer encoder models like ModernBERT (Alshomary et al., 2025b), designing tokenizers specifically for style representations (Wegmann et al., 2025), and pooling not only the last, but several or all, encoder layers might improve performance (Alshomary et al., 2025b).

? Construct interpretable embeddings

An open question is how to learn representations with interpretable dimensions that are still as performant as their uninterpretable counterparts. Such work may benefit from sparse autoencoders—which have recently been shown to automatically learn interpretable features (Huben et al., 2024)—or from combining predefined features with neural training—for example, by training models to classify predefined features (cf. Alkiek et al., 2025).

4 Evaluating style representations

To develop better style representations, we must be able to compare and evaluate them, but no standard currently exists.

4.1 Previous approaches

We divide evaluation approaches according to our definition of style (§2.3), grouping them into **predefined features**, **objects of study**—including authorship verification—and **content-independence**.

On predefined features Learned representations (§3.2) can be evaluated on their sensitivity to predefined features (§3.1). Various studies use probing classifiers (Adi et al., 2017; Köhn, 2015) as well as recurrent/recursive neural networks (Belinkov et al., 2017; Shi et al., 2016) to assess which linguistic features are captured by representations. For example, Alshomary et al., 2025b probe style representations on morphology and syntax. However, probing has limitations, such as uncertainty over how to interpret classifier performance (Belinkov, 2022). Other approaches are sparse, but include studying performance loss on style tasks when re-

moving syntactic and discourse information from texts via shuffling (Zhu and Jurgens, 2021) and evaluating the cosine similarity between texts that include the same predefined features like contraction usage or use of passive voice (Patel et al., 2025; Wegmann and Nguyen, 2021).

On objects of study Representations have also been evaluated for their ability to classify common objects of study (§2), including probing and classifying (i) literary authors (Wang et al., 2023), (ii) book genres (Maharjan et al., 2019), (iii) registers (Alkiek et al., 2025; Wang et al., 2023), and (iv) demographic information of authors like gender or age (Ding et al., 2019; Kang et al., 2019; Kang and Hovy, 2021). Other work examines whether representations of formal/complex texts are similar to other formal/complex texts (Wegmann and Nguyen, 2021). Further, Terreau et al. (2021) use representations to predict an author’s distribution on predefined features.

Authorship attribution Many works evaluate style representations according to their usefulness for authorship attribution or verification tasks (Alkiek et al., 2025; Ding et al., 2019; Maharjan et al., 2019; Patel et al., 2025), including testing whether a representation clusters documents by the same author together (Hay et al., 2020). Datasets and domains like e-mails, blogs, Reddit, Amazon Reviews, Yelp reviews, fanfiction, or shared PAN tasks⁸ from the years 2011–2025 (Argamon and Juola, 2011; Bevendorff et al., 2025a) are commonly used. See Huang et al. (2025) and our GitHub page for a collection of typical datasets. Recently, transcribed spoken domains, such as telephone conversations, interviews, speeches, and podcasts, have also been used (Aggazzotti et al., 2024, 2025b; Tripto et al., 2023). However, without careful preparation, datasets might contain named entities, leakage between train and test sets (Brad et al., 2022; Sawatphol et al., 2024), or spurious correlations with topic (Wegmann et al., 2022), making performance less interpretable. There are promising contributions tackling such issues, like Israeli et al. (2025) and Khan et al. (2021), who provide large sets of authors across different topics on Wikipedia and Reddit, Tripto et al. (2023), who provide speech transcripts across various registers and topics, and Tyo et al. (2022) who design a benchmark across domains for authorship attri-

⁸ <https://pan.webis.de/shared-tasks>

506 bution and verification. However, researchers typi- 555
507 cally use a differing selection of tasks, data, domain 556
508 combinations, or splits, making performance scores 557
509 incomparable across different studies. 558

Content-independence Even though it is debat- 559
510 able whether linguistic style generally excludes con- 560
511 tent information (§2), style representations are com- 561
512 monly tested on “content-independence”. This has 562
513 been evaluated by studying the loss of performance 563
514 on style-related NLP tasks (like authorship verifi- 564
515 cation or attribution) when masking out less fre- 565
516 quent words or “content words” (Stamatatos, 2017; 566
517 Wang et al., 2023; Zhu and Jurgens, 2021) or when 567
518 changing the style of a text with an automatic para- 568
519 phraser (Wang et al., 2023). Other approaches test 569
520 whether style representations are more sensitive 570
521 to style changes than to content changes (Weg- 571
522 mann and Nguyen, 2021; Wegmann et al., 2022), 572
523 whether they can distinguish speakers discussing 573
524 the same conversational topic (Aggazzotti et al., 574
525 2024, 2025b), and whether they perform poorly 575
526 on semantic tasks like topic classification (Wang 576
527 et al., 2023). Generally, few style representations 577
528 reach high scores on content-independence (🔧 578
529 App. Tab. 3) and might benefit from more ex-
530 haustive content disentanglement. 579
531

532 4.2 The future of style evaluation

Increase interpret- and explainability 580
533 The evaluation of learned style representations 581
534 on predefined features is not yet systematic, but 582
535 is promising to pursue, as it can build on rich lit- 583
536 erature in linguistics and stylometry (§2.1, §3.1) 584
537 and can help make learned representations more in- 585
538 terpretable. Further, there is only limited work on 586
539 explaining learned style spaces. Alshomary et al.,
540 2025a pioneer this direction by generating explana-
541 tions on why embeddings cluster certain authors.
542

? Leverage measurement theory 587
543 In the social sciences, measures are commonly 588
544 assessed for *reliability* (do measures return the 589
545 same result with repeated measurement?) and *va-* 590
546 *lidity* (do measures capture the concept of interest?). 591
547 Measurement theory could provide the evaluation 592
548 of style embeddings and the construction of style 593
549 benchmarks with a theoretical framework, highlight 594
550 gaps, and provide inspiration for future methods. 595
551 🔧 See Trochim et al. (2015) for more on measure- 596
552 ment theory. See Fang et al. (2022) for examples of 597
553 how to apply measurement theory to embeddings 598
554

and Bean et al. (2025) for recommendations on how 555
to construct valid benchmarks. We give examples 556
of how measurement theory might be applied for 557
style embeddings and benchmarks in App. §E. 558

Develop standard benchmarks 559
Only a few contributions aim to systematically 560
evaluate representations on linguistic style, leav- 561
ing this area of research behind semantic embed- 562
dings and approaches like MTEB (Enevoldsen et al., 563
2025; Muennighoff et al., 2023). We discuss some 564
notable pioneers: Kang and Hovy (2021) collected 565
the largest dataset to date for style classification; 566
however, several of their classes (e.g., sentiment) 567
would not be considered style in linguistics. Fur- 568
ther, STEL (Wegmann and Nguyen, 2021) is a 569
theory-driven benchmark on single linguistic prop- 570
erties and broader style categories that evaluates 571
representations with cosine similarities—thus not 572
needing training. Neither approach covers a wide 573
range of styles or domains or clearly defines an ob- 574
ject of study (cf. §2.3). Providing an open, easily 575
accessible, high quality, and diverse style benchmark 576
spanning multiple objects of study like registers 577
and authors would be a significant contribution. 578

579 5 What style representations enable


Style representations can make crucial con- 580
tributions to the modern NLP pipeline and 581
to applications of NLP methods. 582

We provide a selection of examples of what style 583
representations can enable. We list a few more in 584
§F, including authorship attribution, bias reduction, 585
reducing spurious correlations in annotations, and 586
improving generalization across styles.

587 5.1 An improved NLP pipeline 588

Curate multi-stage training datasets 589
LLMs are often not robust to stylistic variation 590
(§1). Manipulating and diversifying the style of 591
texts in in-context learning (ICL) examples as well 592
as pre- and post-training datasets—for example, by 593
stratified curation or rephrasing in different styles— 594
has helped output diversity and performance across 595
stylistic variation (Chen et al., 2024b; Lambert et al., 596
2025; Levy et al., 2023; Maini et al., 2024). Curricu- 597
lum learning or multi-stage training found increased 598
success recently (OLMo et al., 2025; Ettinger et al., 599
2025; Allal et al., 2025). We believe that style rep-
resentations can be a crucial tool to monitor the

overall stylistic diversity of a dataset (cf. Nguyen and Ploeger, 2025) and can help select data points for training that increase or decrease stylistic diversity according to a curriculum. Further, they can help rephrase texts in other styles (cf. Maini et al., 2024) using style transfer methods (§5.1) and select datapoints that align with a target style in ICL and (post-)training datasets.

 **Diversify style in evaluation datasets**

Both style and content influence human preference judgments (Cai et al., 2024; Chen et al., 2024a; Singhal et al., 2024; Tianle Li, 2024). However, state-of-the-art performance is often established only on content tasks (mostly NLU and reasoning) using texts with limited stylistic variation (Guo et al., 2025; Truong et al., 2025). This might obfuscate a model’s ability to generalize to other domains or understand and generate diverse or preferred styles.⁹ Instead, benchmarks could be composed not only based on what they test, but also based on whether their datasets or tasks cover different or expected regions of the style embedding space.

5.2 Various other applications

Generating in specific styles Representations of style can help generate text in a specific style, or adapt to different domains (Horvitz et al., 2024a,b; Liu et al., 2023; Zhang et al., 2023a). Such style steering approaches can enable accessibility in language generation (Anschütz et al., 2025; Cao et al., 2020; Surya et al., 2019)—for example, by simplifying a text for a child or summarizing a text for a non-expert. The style of generated texts is often evaluated by comparing their style representations to those of a target style (Chim et al., 2025; Horvitz et al., 2024a; Jangra et al., 2025; Liu et al., 2023).

Personalization Interest in personalized model responses has grown recently (Zhang et al., 2025b; Liu et al., 2025). Style plays a crucial role in personalization (Zhang et al., 2025b; Liu et al., 2025), and style representations could be used to recognize the style of humans, infer their preferences, and adapt generated responses to them (Zhang et al., 2025a).

Machine text detection There is a growing concern about the misuse of LLMs, including disinformation, spam, and plagiarism. Recent work (Bevendorff et al., 2025b; Elkhatat et al., 2023; Gehrmann et al., 2019; Kumarage et al., 2023; Sun et al., 2025;

⁹For example, textbooks might not be all you need (cf. Li et al., 2023) for perplexity across registers (Maini, 2023).

Uchendu et al., 2020) shows that LLMs exhibit idiosyncrasies that distinguish their writing from human writing. Detectors that use style embeddings have been effective in in-domain and cross-domain settings (Kim et al., 2025; Rivera Soto et al., 2023).

Privacy On the flip side of attribution and detection is the task of obfuscating someone’s identity.¹⁰ Style representations can help determine if text that has been anonymized, such as via paraphrasing, sufficiently removes someone’s style and protects their privacy (Aggazzotti et al., 2025a; Alperin et al., 2025; Bao and Carpuat, 2024; Shokri et al., 2025).

 **Push style representations as a foundational method for NLP and other fields**

Just as semantic embeddings have become foundational, style representations could also be foundational across fields. Next to the mentioned uses, they could help retrieve documents with a (dis-)similar style to a search query (Cao, 2025), track style shift in dialogue in sociolinguistics (Nguyen, 2025), or analyze literary text in the digital humanities (Hicke and Mimno, 2025), with current embeddings already seeing significant adoption.¹¹

6 Conclusion

With this paper, we hope to have demonstrated the potential of style representations for the NLP community. We call on researchers to use clearer definitions of style, to more explicitly disentangle evaluation and training approaches, and to develop evaluation methods into a standard. We end by noting that style has unique properties that may require unique considerations and methodologies. Among these, the style of a text is inherently relative. For example, it might be clearer and more relevant to judge if a text (e.g., *How are you?*) is more formal than another (e.g., *What’s up?*) rather than if it is formal in isolation; consider also App. Fig. 3 and Irvine (2001). This relativity may require new solutions in training and evaluating representations—for example, curating training data with hard positives and negatives positioned in relation to each other, or testing whether representations correctly rank sentences along a stylistic dimension.

¹⁰For example, see the PAN Author Masking series at <https://pan.webis.de/shared-tasks.html#author-masking>.

¹¹<https://huggingface.co/AnnaWegmann/Style-Embedding> reached 200k downloads in October 2025.

689 Limitations

690 Consider style in modalities other than 691 text.

692 Many of the examples and citations throughout
693 this paper refer to text-based style since the lim-
694 ited style research in NLP has focused on written
695 language, but linguistic style also manifests, and
696 is perhaps better studied, in other modalities like
697 speech (e.g., tone of voice), gestures, and vision
698 (e.g., image generation). We leave considerations
699 for representing style in other modalities for future
work.

700 Give more attention to style in languages 701 other than English.

702 The bulk of the work we discuss considers style
703 in English. For example, we mainly discuss def-
704 initions of style considered by American schol-
705 ars (cf. §2.1), and we discuss predefined features
706 mainly for English (cf. §3.1)—for instance, “g-
707 dropping” is an English-specific marker. Different
708 scripts and languages will usually need different
709 predefined features and have a different history
710 regarding style definitions and sociolinguistic re-
711 search (see also Ball et al., 2023). However, our
712 discussed approaches to automatically learn and
713 evaluate representations should largely transcend
714 languages and scripts as long as architectural com-
715 ponents (e.g., tokenizers), evaluation datasets, and
716 predefined features are adapted for optimal perfor-
717 mance. We believe that developing style represen-
718 tations for languages other than English is a crucial
719 future step and call on the community to continue
720 pioneering work like Kim et al. (2025) and Qiu et al.
(2025).

721 Why not use a different term instead of style?

722 *...the extremely broad and ambiguous ref-*
723 *erence of the term [style] in everyday use*
724 *has not made its status as a technical lin-*
725 *guistic term very appealing.*

726 — David Crystal

727 Scholars, such as Crystal (2011), have argued
728 against using the term style at all due to its increas-
729 ingly vague and colloquial use. Instead, researchers
730 have opted to describe the specific phenomenon
731 they are interested in (e.g., syntactic variation) and
732 use less over-defined terms (e.g., language varia-
733 tion). While that can be helpful in some cases, we
734 argue that using the term style is still worthwhile
735 because (i) the term is used regularly in NLP (with

200 publications in the ACL Anthology mentioning
“style” in the title or abstract in 2024) highlighting
the interest in the term; (ii) style seems to provide
a more concise and intuitive label than alternatives
like “distinctive patterns in the used language vari-
eties” or “systematic variation in textual features”;
and (iii) the term style can draw from decades of
theoretical foundation in stylometry and sociolin-
guistics.

Style is a concept used in many fields. Why fo-
cus on the ones discussed in the paper? Next
to NLP, we focus on definitions and concepts of
style used in sociolinguistics, linguistics, stylome-
try, forensic linguistics, and corpus linguistics (§2,
see an overview of the fields in §C). To the best
of our knowledge, these are the most active areas
already using, or intuitive areas that could profit
from using, computational methods for analyzing
style. Further, we believe that sociolinguistics is
particularly relevant to consider, as its study of
the interaction between language and society has
unique potential to inform NLP methods (Nguyen,
2025), especially as NLP models are increasingly
used within, and have growing impact on, society.

760 Ethical considerations

761 Style modeling is closely related to *author profiling*
762 (cf. §4)—the task of recovering author character-
763 istics based on the text they wrote (Nguyen et al.,
764 2013; Rangel et al., 2013). Note that author pro-
765 filing can be useful for improving performance on
766 some NLP tasks (Hovy, 2015); however, identifi-
767 ing an author’s gender, age, personality type, etc.
768 has increasingly been criticized for bias and pri-
769 vacy concerns (Brennan et al., 2012; Elazar and
770 Goldberg, 2018; Li et al., 2018; Lison et al., 2021).

771 Integrating more language diversity, and with it
772 social factors, into NLP is a double-edged sword:
773 There are clear advantages to integrating more di-
774 versity into NLP models and, specifically, repre-
775 senting minorities to increase the fairness and rep-
776 resentativeness of NLP models (Bird and Yibarbuk,
777 2024; Grieve et al., 2025; Hovy and Yang, 2021;
778 Markl et al., 2024); however, making NLP models
779 more sensitive to social factors could also make
780 them a threat to privacy across social groups. The
781 performance of machine learning approaches on
782 tasks like author profiling could increase, result-
783 ing in a large potential for misuse, such as the following
784 examples: (1) Author profiles could be used to iden-
785 tify and profile individuals or political dissenters

(Hovy and Spruit, 2016); (2) Author profiling could be used for predatory ad targeting, which might show gambling ads to vulnerable groups or not show job postings to certain social groups (Dudy et al., 2021); and (3) Author profiles could lead to data leakage, such as making health conditions recoverable for insurance companies that might increase their rates for certain individuals (Dudy et al., 2021).

This conflict between privacy and fairness has been described as one of the “dual-use problems” in NLP by Hovy and Spruit (2016). We aim to improve fairness without compromising individual privacy and safety but acknowledge that progress in one might sometimes come at the expense of the other. 🚩 Therefore, we encourage researchers in the NLP community to engage with the dual-use problem more actively and work on techniques to make the design of language models more sensitive to human values, as suggested in Dudy et al. (2021), ideally without actively working on approaches to make sensitive data recoverable from texts. We further encourage researchers to actively anonymize datasets used for modeling and the evaluation of style representations.

We confirm that we have read and abide by the ACL Code of Ethics. Besides those mentioned, we do not foresee immediate risks of our work.

References

Ahmed Abbasi and Hsinchun Chen. 2008. [Writeprints: A stylometric approach to identity-level identification and similarity detection in cyberspace](#). *ACM Transactions on Information Systems (TOIS)*, 26(2):1–29.

Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi, and Yoav Goldberg. 2017. [Fine-grained analysis of sentence embeddings using auxiliary prediction tasks](#). In *International Conference on Learning Representations (ICLR)*.

Cristina Aggazzotti, Nicholas Andrews, and Elizabeth Allyn Smith. 2024. [Can authorship attribution models distinguish speakers in speech transcripts?](#) *Transactions of the Association for Computational Linguistics*, 12:875–891.

Cristina Aggazzotti, Ashi Garg, Zexin Cai, and Nicholas Andrews. 2025a. [Content anonymization for privacy in long-form audio](#). *arXiv preprint*. ArXiv:2510.12780 [cs].

Cristina Aggazzotti and Elizabeth Allyn Smith. 2025. [A stylometric analysis of speaker attribution from speech transcripts](#). *Preprint*, arXiv:2512.13667.

Cristina Aggazzotti, Matthew Wiesner, Elizabeth Allyn Smith, and Nicholas Andrews. 2025b. [The impact of automatic speech transcription on speaker attribution](#). *Transactions of the Association for Computational Linguistics*, in press.

Kenan Alkiek, Anna Wegmann, Jian Zhu, and David Jurgens. 2025. [Neurobiber: Fast and interpretable stylistic feature extraction](#). *arXiv preprint*. ArXiv:2502.18590 [cs].

Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martin Blazquez, Guilherme Penedo, Lewis Tunstall, Andrés Marafioti, Agustín Piqueres Lajarin, Hynek Kydlíček, Vaibhav Srivastav, Joshua Lochner, Caleb Fahlgren, Xuan Son Nguyen, Ben Burtenshaw, Clémentine Fourrier, Haojun Zhao, Hugo Larcher, Mathieu Morlon, Cyril Zakka, and 3 others. 2025. [SmolLM2: When smol goes big — Data-centric training of a fully open small language model](#). In *Second Conference on Language Modeling (COLM)*.

Kenneth Alperin, Rohan Leekha, Adaku Uchendu, Trang Nguyen, Srilakshmi Medarametla, Carlos Levya Capote, Seth Aycock, and Charlie Dagli. 2025. [Masks and mimicry: Strategic obfuscation and impersonation attacks on authorship verification](#). In *Proceedings of the 5th International Conference on Natural Language Processing for Digital Humanities*, pages 102–116, Albuquerque, USA. Association for Computational Linguistics.

Milad Alshomary, Narutatsu Ri, Marianna Apidianaki, Ajay Patel, Smaranda Muresan, and Kathleen McKeown. 2025a. [Latent space interpretation for stylistic analysis and explainable authorship attribution](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 1124–1135, Abu Dhabi, UAE. Association for Computational Linguistics.

Milad Alshomary, Nikhil Reddy Varimalla, Vishal Anand, Smaranda Muresan, and Kathleen McKeown. 2025b. [Layered insights: Generalizable analysis of human authorial style by leveraging all transformer layers](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 10290–10303, Suzhou, China. Association for Computational Linguistics.

Malik Altkrori, Jackie Chi Kit Cheung, and Benjamin CM Fung. 2021. [The topic confusion task: A novel evaluation scenario for authorship attribution](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4242–4256.

Nicholas Andrews and Marcus Bishop. 2019. [Learning invariant representations of social media users](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1684–1695, Hong Kong, China. Association for Computational Linguistics.

893	Miriam Anshütz, Anastasiya Damaratskaya,	Batzner, Negar Foroutan, Chris Schmitz, Karolina	950
894	Chaeun Joy Lee, Arthur Schmalz, Edoardo	Korgul, Hunar Batra, Oishi Deb, Emma Beharry, Cor-	951
895	Mosca, and Georg Groh. 2025. (Dis)improved?! How simplified language affects large language	nelius Emde, Thomas Foster, Anna Gausen, Maria	952
896	model performance across languages. In <i>Pro-</i>	Grandury, Simeng Han, Valentin Hofmann, Lujain	953
897	<i>ceedings of the Fourth Workshop on Generation,</i>	Ibrahim, and 23 others. 2025. Measuring what mat-	954
898	<i>Evaluation and Metrics (GEM²)</i> , pages 847–861,	ters: Construct validity in large language model	955
899	Vienna, Austria and virtual meeting. Association for	benchmarks. In <i>The Thirty-ninth Annual Confer-</i>	956
900	Computational Linguistics.	<i>ence on Neural Information Processing Systems</i>	957
901		<i>(NeurIPS)</i> .	958
902	Ehsan Arabnezhad, Massimo La Morgia, Alessandro	Yonatan Belinkov. 2022. Probing classifiers: Promises,	959
903	Mei, Eugenio Nerio Nemmi, and Julinda Stefa. 2020.	shortcomings, and advances. <i>Computational Linguis-</i>	960
904	A light in the dark web: Linking dark web aliases to	<i>tics</i> , 48(1):207–219.	961
905	real internet identities. In <i>Proceedings of the 40th</i>	Yonatan Belinkov, Nadir Durrani, Fahim Dalvi, Has-	962
906	<i>International Conference on Distributed Computing</i>	san Sajjad, and James Glass. 2017. What do neural	963
907	<i>Systems (ICDCS)</i> , pages 311–321, Singapore, Singa-	machine translation models learn about morphology?	964
908	pore. Institute of Electrical and Electronics Engineers.	In <i>Proceedings of the 55th Annual Meeting of the As-</i>	965
909	Shlomo Argamon. 2018. Computational forensic au-	<i>sociation for Computational Linguistics (Volume 1:</i>	966
910	thorship analysis: Promises and pitfalls. <i>Language</i>	<i>Long Papers)</i> , pages 861–872, Vancouver, Canada.	967
911	<i>and Law/Linguagem e Direito</i> , 5(2):7–37.	Association for Computational Linguistics.	968
912	Shlomo Argamon and Patrick Juola. 2011. Overview of	Allan Bell. 1984. Language style as audience design.	969
913	the international authorship identification competition	<i>Language in Society</i> , 13(2):145–204.	970
914	at PAN-2011. In <i>Notebook Papers of CLEF 2011</i>	Allan Bell. 2014. <i>The Guidebook to Sociolinguistics.</i>	971
915	<i>Labs and Workshops</i> , Amsterdam, Netherlands.	John Wiley & Sons, Chichester, UK.	972
916	Mikel Artetxe, Shruti Bhosale, Naman Goyal, Todor	Janek Bevendorff, Daryna Dementieva, Maik Fröbe,	973
917	Mihaylov, Myle Ott, Sam Shleifer, Xi Victoria Lin,	Bela Gipp, André Greiner-Petter, Jussi Karlgren,	974
918	Jingfei Du, Srinivasan Iyer, Ramakanth Pasunuru,	Maximilian Mayerl, Preslav Nakov, Alexander	975
919	Giridharan Anantharaman, Xian Li, Shuohui Chen,	Panchenko, Martin Potthast, Artem Shelmanov, Efs-	976
920	Halil Akin, Mandeep Baines, Louis Martin, Xing	tathios Stamatatos, Benno Stein, Yuxia Wang, Matti	977
921	Zhou, Punit Singh Koura, Brian O’Horo, and 5 oth-	Wiegmann, and Eva Zangerle. 2025a. Overview of	978
922	ers. 2022. Efficient large scale language modeling	PAN 2025: Generative AI detection, multilingual	979
923	with mixtures of experts. In <i>Proceedings of the 2022</i>	text detoxification, multi-author writing style analy-	980
924	<i>Conference on Empirical Methods in Natural Lan-</i>	sis, and generative plagiarism detection. In <i>Advances</i>	981
925	<i>guage Processing</i> , pages 11699–11732.	<i>in Information Retrieval</i> , pages 434–441. Springer,	982
926	Martin J. Ball, Rajend Mesthrie, and Chiara Meluzzi.	Cham.	983
927	2023. <i>The Routledge Handbook of Sociolinguistics</i>	Janek Bevendorff, Matti Wiegmann, Emmelie Richter,	984
928	<i>Around the World</i> , 2nd edition. Routledge, London,	Martin Potthast, and Benno Stein. 2025b. The two	985
929	UK.	paradigms of LLM detection: Authorship attribution	986
930	Akshina Banerjee and Oleg Urminsky. 2025. The lan-	vs. authorship verification. In <i>Findings of the Asso-</i>	987
931	guage that drives engagement: A systematic large-	<i>ciation for Computational Linguistics: ACL 2025,</i>	988
932	scale analysis of headline experiments. <i>Marketing</i>	pages 3762–3787, Vienna, Austria. Association for	989
933	<i>Science</i> , 44(3):566–592.	Computational Linguistics.	990
934	Calvin Bao and Marine Carpuat. 2024. Keep it Private:	Douglas Biber. 1988. <i>Variation across Speech and Writ-</i>	991
935	Unsupervised privatization of online text. In <i>Proceed-</i>	<i>ing.</i> Cambridge University Press, Cambridge, UK.	992
936	<i>ings of the 2024 Conference of the North American</i>	Douglas Biber and Susan Conrad. 2019. <i>Register,</i>	993
937	<i>Chapter of the Association for Computational Lin-</i>	<i>Genre, and Style</i> , 2nd edition. Cambridge University	994
938	<i>guistics: Human Language Technologies (Volume</i>	Press, Cambridge, UK.	995
939	<i>1: Long Papers)</i> , pages 8678–8693, Mexico City,	Steven Bird, Ewan Klein, and Edward Loper. 2019. <i>Nat-</i>	996
940	Mexico. Association for Computational Linguistics.	<i>ural Language Processing with Python.</i> O’Reilly	997
941	Yu Bao, Hao Zhou, Shujian Huang, Lei Li, Lili Mou,	Media.	998
942	Olga Vechtomova, Xin-yu Dai, and Jiajun Chen.	Steven Bird and Dean Yibarbuk. 2024. Centering the	999
943	2019. Generating sentences from disentangled syn-	speech community. In <i>Proceedings of the 18th Con-</i>	1000
944	tactic and semantic spaces. In <i>Proceedings of the</i>	<i>ference of the European Chapter of the Association</i>	1001
945	<i>57th Annual Meeting of the Association for Computa-</i>	<i>for Computational Linguistics (Volume 1: Long Pa-</i>	1002
946	<i>tional Linguistics</i> , pages 6008–6019, Florence, Italy.	<i>pers)</i> , pages 826–839, St. Julian’s, Malta. Association	1003
947	Association for Computational Linguistics.	for Computational Linguistics.	1004
948	Andrew M. Bean, Ryan Othniel Kearns, Angelika Ro-		
949	manou, Franziska Sofia Hafner, Harry Mayne, Jan		

1224	alignment benchmarking . In <i>The Thirteenth International Conference on Learning Representations (ICLR)</i> .	<i>Web</i> , pages 303–322. Springer, Dordrecht, the Netherlands.	1278
1225			1279
1226			
1227	Eve Fleisig, Genevieve Smith, Madeline Bossi, Ishita Rustagi, Xavier Yin, and Dan Klein. 2024. Linguistic bias in ChatGPT: Language models reinforce dialect discrimination . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 13541–13564, Miami, Florida, USA. Association for Computational Linguistics.	Yanzhu Guo, Guokan Shang, and Chloé Clavel. 2025. Benchmarking linguistic diversity of large language models . <i>arXiv preprint</i> . ArXiv:2412.10271 [cs].	1280
1228			1281
1229			1282
1230		Yanzhu Guo, Guokan Shang, Michalis Vazirgiannis, and Chloé Clavel. 2024. The curious decline of linguistic diversity: Training language models on synthetic text . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 3589–3604, Mexico City, Mexico. Association for Computational Linguistics.	1283
1231			1284
1232			1285
1233			1286
1234	Albert Gatt and Emiel Kraemer. 2018. Survey of the state of the art in natural language generation: Core tasks, applications and evaluation . <i>Journal of Artificial Intelligence Research</i> , 61:65–170.		1287
1235			1288
1236			1289
1237		Suchin Gururangan, Swabha Swayamdipta, Omer Levy, Roy Schwartz, Samuel Bowman, and Noah A. Smith. 2018. Annotation artifacts in natural language inference data . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 107–112, New Orleans, Louisiana. Association for Computational Linguistics.	1290
1238	Sebastian Gehrmann, Hendrik Strobelt, and Alexander Rush. 2019. GLTR: Statistical detection and visualization of generated text . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> , pages 111–116, Florence, Italy. Association for Computational Linguistics.		1291
1239			1292
1240			1293
1241			1294
1242			1295
1243			1296
1244			1297
1245	Howard Giles, Nikolas Coupland, and Justine Coupland. 1991. Accommodation theory: Communication, context, and consequence . <i>Contexts of accommodation: Developments in applied sociolinguistics</i> , 1:1–68.	Julien Hay, Bich-Lien Doan, Fabrice Popineau, and Ouassim Ait Elhara. 2020. Representation learning of writing style . In <i>Proceedings of the Sixth Workshop on Noisy User-generated Text (W-NUT 2020)</i> , pages 232–243, Online. Association for Computational Linguistics.	1299
1246			1300
1247			1301
1248			1302
1249	Howard Giles and Peter F. Powesland. 1975. <i>Speech Style and Social Evaluation</i> . Academic Press, London, UK.		1303
1250			1304
1251		Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset . In <i>Thirty-fifth Conference on Neural Information Processing Systems (NeurIPS)</i> .	1305
1252	Yoav Goldberg. 2019. Assessing BERT’s syntactic abilities . <i>arXiv preprint</i> . ArXiv:1901.05287 [cs].		1306
1253			1307
1254	Arthur C. Graesser, Danielle S. McNamara, Max M. Louwerse, and Zhiqiang Cai. 2004. Coh-Matrix: Analysis of text on cohesion and language . <i>Behavior Research Methods, Instruments, & Computers</i> , 36(2):193–202.		1308
1255			1309
1256		Rebecca M. M. Hicke and David Mimno. 2025. Looking for the inner music: Probing LLMs’ understanding of literary style . <i>Computational Humanities Research</i> , 1:e3. Publisher: Cambridge University Press.	1311
1257			1312
1258			1313
1259	Tim Grant. 2022. <i>The Idea of Progress in Forensic Authorship Analysis</i> . Elements in Forensic Linguistics. Cambridge University Press.		1314
1260		Valentin Hofmann, Pratyusha Ria Kalluri, Dan Jurafsky, and Sharese King. 2024. AI generates covertly racist decisions about people based on their dialect . <i>Nature</i> , 633:147–154.	1315
1261			1316
1262	Jack Grieve. 2007. Quantitative authorship attribution: An evaluation of techniques . <i>Literary and Linguistic Computing</i> , 22(3):251–270.		1317
1263		Nicole Holliday. 2021. Intonation and referee design phenomena in the narrative speech of Black/biracial men . <i>Journal of English Linguistics</i> , 49(3):283–304.	1319
1264			1320
1265	Jack Grieve. 2023. Register variation explains stylistometric authorship analysis . <i>Corpus Linguistics and Linguistic Theory</i> , 19(1):47–77.		1321
1266		David I. Holmes. 1985. The analysis of literary style—A review . <i>Journal of the Royal Statistical Society: Series A (General)</i> , 148(4):328–341.	1322
1267			1323
1268	Jack Grieve, Sara Bartl, Matteo Fuoli, Jason Grafmiller, Weihang Huang, Alejandro Jawerbaum, Akira Murakami, Marcus Perlman, Dana Roemling, and Bodo Winter. 2025. The sociolinguistic foundations of language modeling . <i>Frontiers in Artificial Intelligence</i> , 7:1472411.		1324
1269		Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spaCy: Industrial-strength natural language processing in Python .	1325
1270			1326
1271		Zachary Horvitz, Ajay Patel, Chris Callison-Burch, Zhou Yu, and Kathleen McKeown. 2024a. ParaGuide: Guided diffusion paraphrasers for plug-and-play textual style transfer . In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , pages	1327
1272			1328
1273			1329
1274	Jack Grieve, Douglas Biber, Eric Friginal, and Tatiana Nekrasova. 2011. Variation among blogs: A multi-dimensional analysis . In Alexander Mehler, Serge Sharoff, and Marina Santini, editors, <i>Genres on the</i>		1330
1275			1331
1276			1332
1277			

1333	18216–18224, Vancouver, Canada. Association for the Advancement of Artificial Intelligence.	1389
1334		1390
1335	Zachary Horvitz, Ajay Patel, Kanishk Singh, Chris Callison-Burch, Kathleen McKeown, and Zhou Yu. 2024b. TinyStyler: Efficient few-shot text style transfer with authorship embeddings . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 13376–13390, Miami, Florida, USA. Association for Computational Linguistics.	1391
1336		1392
1337		1393
1338		1394
1339		1395
1340		1396
1341		1397
1342	John Houvardas and Efstathios Stamatatos. 2006. N-gram feature selection for authorship identification . In <i>Proceedings of the 12th International Conference on Artificial Intelligence: Methodology, Systems, and Applications</i> , AIMS’06, page 77–86, Berlin, Germany. Springer.	1398
1343		1399
1344		1400
1345		1401
1346		1402
1347		1403
1348	Dirk Hovy. 2015. Demographic factors improve classification performance . In <i>Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 752–762, Beijing, China. Association for Computational Linguistics.	1404
1349		1405
1350		1406
1351		1407
1352		1408
1353		1409
1354		1410
1355	Dirk Hovy, Federico Bianchi, and Tommaso Fornaciari. 2020. “You sound just like your father” commercial machine translation systems include stylistic biases . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1686–1690, Online. Association for Computational Linguistics.	1411
1356		1412
1357		1413
1358		1414
1359		1415
1360		1416
1361		1417
1362	Dirk Hovy and Shannon L. Spruit. 2016. The social impact of natural language processing . In <i>Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 591–598, Berlin, Germany. Association for Computational Linguistics.	1418
1363		1419
1364		1420
1365		1421
1366		1422
1367		1423
1368	Dirk Hovy and Diyi Yang. 2021. The importance of modeling social factors of language: Theory and practice . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 588–602, Online. Association for Computational Linguistics.	1424
1369		1425
1370		1426
1371		1427
1372		1428
1373		1429
1374		1430
1375	Baixiang Huang, Canyu Chen, and Kai Shu. 2025. Authorship attribution in the era of LLMs: Problems, methodologies, and challenges . <i>ACM SIGKDD Explorations Newsletter</i> , 26(2):21–43.	1431
1376		1432
1377		1433
1378		1434
1379	Robert Huben, Hoagy Cunningham, Logan Riggs Smith, Aidan Ewart, and Lee Sharkey. 2024. Sparse autoencoders find highly interpretable features in language models . In <i>The Twelfth International Conference on Learning Representations</i> .	1435
1380		1436
1381		1437
1382		1438
1383		1439
1384	Judith T. Irvine. 2001. “Style” as distinctiveness: the culture and ideology of linguistic differentiation . In Penelope Eckert and John R. Rickford, editors, <i>Style and Sociolinguistic Variation</i> , pages 21–43. Cambridge University Press, Cambridge, UK.	1440
1385		1441
1386		1442
1387		1443
1388		1444
		1445
	Abraham Israeli, Shuai Liu, Jonathan May, and David Jurgens. 2025. The Million Authors corpus: A cross-lingual and cross-domain Wikipedia dataset for authorship verification . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 2597–26017, Vienna, Austria. Association for Computational Linguistics.	1446
		1447
	Somayeh Jafaritazehjani, GwénoLé Lecorvé, Damien Lolive, and John Kelleher. 2020. Style versus content: A distinction without a (learnable) difference? In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 2169–2180, Barcelona, Spain (Online). International Committee on Computational Linguistics.	1448
		1449
	Anubhav Jangra, Bahareh Sarrafzadeh, Adrian de Wynter, Silviu Cucerzan, and Sujay Kumar Jauhar. 2025. Evaluating style-personalized text generation: Challenges and directions . <i>arXiv preprint</i> . ArXiv:2508.06374 [cs].	1450
		1451
	Harsh Jhamtani, Varun Gangal, Eduard Hovy, and Eric Nyberg. 2017. Shakespeareizing modern language using copy-enriched sequence to sequence models . In <i>Proceedings of the Workshop on Stylistic Variation</i> , pages 10–19, Copenhagen, Denmark. Association for Computational Linguistics.	1452
		1453
	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7B . <i>arXiv preprint</i> . ArXiv:2310.06825 [cs].	1454
		1455
	Di Jin, Zhijing Jin, Zhiting Hu, Olga Vechtomova, and Rada Mihalcea. 2022. Deep learning for text style transfer: A survey . <i>Computational Linguistics</i> , 48(1):155–205.	1456
		1457
	Vineet John, Lili Mou, Hareesh Bahuleyan, and Olga Vechtomova. 2019. Disentangled representation learning for non-parallel text style transfer . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 424–434, Florence, Italy. Association for Computational Linguistics.	1458
		1459
	Patrick Juola. 2006. Authorship attribution . <i>Foundations and Trends in Information Retrieval</i> , 1(3):233–334.	1460
		1461
	Patrick Juola, John Noecker Jr., Mike Ryan, and Sandy Speer. 2009. JGAAP 4.0—A revised authorship attribution tool. <i>Proceedings of Digital Humanities</i> .	1462
		1463
	Dongyeop Kang, Varun Gangal, and Eduard Hovy. 2019. (male, bachelor) and (female, Ph.D) have different connotations: Parallely annotated stylistic language dataset with multiple personas . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing</i>	1464
		1465

1446		(EMNLP-IJCNLP), pages 1696–1706, Hong Kong, China. Association for Computational Linguistics.	
1447			
1448	Dongyeop Kang and Eduard Hovy. 2021. Style is NOT a single variable: Case studies for cross-stylistic language understanding . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 2376–2387, Online. Association for Computational Linguistics.		
1449			
1450			
1451			
1452			
1453			
1454			
1455			
1456	Mike Kestemont. 2014. Function words in authorship attribution. From black magic to theory? In <i>Proceedings of the 3rd Workshop on Computational Linguistics for Literature (CLFL)</i> , pages 59–66, Gothenburg, Sweden. Association for Computational Linguistics.		
1457			
1458			
1459			
1460			
1461	Salar Khaleghzadegan, Michael Rosen, Anne Links, Alya Ahmad, Molly Kilcullen, Emily Boss, Mary Catherine Beach, and Somnath Saha. 2024. Validating computer-generated measures of linguistic style matching and accommodation in patient-clinician communication . <i>Patient Education and Counseling</i> , 119:108074.		
1462			
1463			
1464			
1465			
1466			
1467			
1468	Aleem Khan, Elizabeth Fleming, Noah Schofield, Marcus Bishop, and Nicholas Andrews. 2021. A deep metric learning approach to account linking . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5275–5287, Online. Association for Computational Linguistics.		
1469			
1470			
1471			
1472			
1473			
1474			
1475			
1476	Aleem Khan, Andrew Wang, Sophia Hager, and Nicholas Andrews. 2024. Learning to generate text in arbitrary writing styles . <i>arXiv preprint</i> . ArXiv:2312.17242 [cs].		
1477			
1478			
1479			
1480	Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. 2020. Supervised contrastive learning . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 18661–18673. Curran Associates, Inc.		
1481			
1482			
1483			
1484			
1485			
1486	Junghwan Kim, Haotian Zhang, and David Jurgens. 2025. Leveraging multilingual training for authorship representation: Enhancing generalization across languages and domains . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 34855–34880, Suzhou, China. Association for Computational Linguistics.		
1487			
1488			
1489			
1490			
1491			
1492			
1493	Hannes Kniffka. 2007. <i>Working in Language and Law: A German Perspective</i> . Palgrave Macmillan UK.		
1494			
1495	Arne Köhn. 2015. What’s in an embedding? Analyzing word embeddings through multilingual evaluation . In <i>Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing</i> , pages 2067–2073, Lisbon, Portugal. Association for Computational Linguistics.		
1496			
1497			
1498			
1499			
1500			
	Moshe Koppel, Shlomo Argamon, and Anat Rachel Shmuni. 2002. Automatically categorizing written texts by author gender . <i>Literary and Linguistic Computing</i> , 17(4):401–412.		1501 1502 1503 1504
	Tore Kristiansen. 2024. Social variation in germanic . <i>Oxford Research Encyclopedia of Linguistics</i> .		1505 1506
	Tharindu Kumarage, Joshua Garland, Amrita Bhat-tacharjee, Kirill Trapeznikov, Scott Ruston, and Huan Liu. 2023. Stylometric detection of AI-generated text in Twitter timelines . <i>arXiv preprint</i> . ArXiv:2303.03697 [cs].		1507 1508 1509 1510 1511
	William Labov. 1972. <i>Sociolinguistic Patterns</i> . University of Pennsylvania Press, Philadelphia, USA.		1512 1513
	William Labov. 2006. <i>The Social Stratification of English in New York City</i> , 2nd edition. Cambridge University Press, Cambridge, UK.		1514 1515 1516
	Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James Validad Miranda, Alisa Liu, Nouha Dziri, Xinxu Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Christopher Wilhelm, Luca Soldaini, and 4 others. 2025. Tulu 3: Pushing frontiers in open language model post-training . In <i>Second Conference on Language Modeling (COLM)</i> .		1517 1518 1519 1520 1521 1522 1523 1524 1525
	Beatriz R. Lavandera. 1978. Where does the sociolinguistic variable stop? <i>Language in Society</i> , 7(2):171–182.		1526 1527 1528
	Bruce W. Lee and Jason Lee. 2023. LFTK: Handcrafted features in computational linguistics . In <i>Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)</i> , pages 1–19, Toronto, Canada. Association for Computational Linguistics.		1529 1530 1531 1532 1533 1534
	Itay Levy, Ben Bogin, and Jonathan Berant. 2023. Diverse demonstrations improve in-context compositional generalization . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1401–1422, Toronto, Canada. Association for Computational Linguistics.		1535 1536 1537 1538 1539 1540 1541
	Jinfeng Li, Shouling Ji, Tianyu Du, Bo Li, and Ting Wang. 2019. TextBugger: Generating adversarial text against real-world applications . In <i>Proceedings 2019 Network and Distributed System Security Symposium</i> , San Diego, CA. Internet Society.		1542 1543 1544 1545 1546
	Yitong Li, Timothy Baldwin, and Trevor Cohn. 2018. Towards robust and privacy-preserving text representations . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 25–30, Melbourne, Australia. Association for Computational Linguistics.		1547 1548 1549 1550 1551 1552
	Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat Lee. 2023. Textbooks are all you need II: Phi-1.5 technical report . <i>arXiv preprint</i> . ArXiv:2309.05463 [cs].		1553 1554 1555 1556

1557	Weixin Liang, Mert Yuksekgonul, Yining Mao, Eric Wu, and James Zou. 2023. GPT detectors are biased against non-native English writers . <i>Patterns</i> , 4(7):100779.	1614
1558		1615
1559		1616
1560		1617
1561	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. 2023. Let’s verify step by step . In <i>The Twelfth International Conference on Learning Representations (ICLR)</i> .	1618
1562		1619
1563		1620
1564		1621
1565		1622
1566		1623
1567	Philip Lippmann and Jie Yang. 2025. Style over substance: Distilled language models reason via stylistic replication . In <i>Second Conference on Language Modeling (COLM)</i> .	1624
1568		1625
1569		1626
1570		1627
1571		1628
1572	Pierre Lison, Ildikó Pilán, David Sanchez, Montserrat Batet, and Lilja Øvrelid. 2021. Anonymisation models for text data: State of the art, challenges and future directions . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 4188–4203, Online. Association for Computational Linguistics.	1629
1573		1630
1574		1631
1575		1632
1576		1633
1577		1634
1578		1635
1579		
1580	Frederick Liu, Terry Huang, Shihang Lyu, Siamak Shakeri, Hongkun Yu, and Jing Li. 2022. Enct5: A framework for fine-tuning t5 as non-autoregressive models . <i>arXiv preprint</i> . ArXiv:2110.08426 [cs].	1636
1581		1637
1582		1638
1583		1639
1584	Jiahong Liu, Zexuan Qiu, Zhongyang Li, Quanyu Dai, Wenhao Yu, Jieming Zhu, Minda Hu, Menglin Yang, Tat-Seng Chua, and Irwin King. 2025. A survey of personalized large language models: Progress and future directions . <i>arXiv preprint</i> . ArXiv:2502.11528 [cs].	1640
1585		1641
1586		1642
1587		1643
1588		1644
1589		1645
1590	Shuai Liu, Hyundong Cho, Marjorie Freedman, Xuezhe Ma, and Jonathan May. 2023. RECAP: Retrieval-Enhanced Context-Aware Prefix encoder for personalized dialogue response generation . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8404–8419, Toronto, Canada. Association for Computational Linguistics.	1646
1591		1647
1592		1648
1593		1649
1594		1650
1595		1651
1596		1652
1597		1653
1598		1654
1599		1655
1600	Stephan Ludwig, Ko de Ruyter, Max Friedman, Elisabeth Constantin Brüggem, Martin Wetzels, and Gerard Pfann. 2013. More than words: The influence of affective content and linguistic style matches in online reviews on conversion rates . <i>Journal of Marketing</i> , 77(1):87–103.	1656
1601		1657
1602		1658
1603		1659
1604	Suraj Maharjan, Deepthi Mave, Prasha Shrestha, Manuel Montes, Fabio A. González, and Thamar Solorio. 2019. Jointly learning author and annotated character n-gram embeddings: A case study in literary text . In <i>Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2019)</i> , pages 684–692, Varna, Bulgaria. INCOMA Ltd.	1660
1605		1661
1606		1662
1607		1663
1608		1664
1609		1665
1610		1666
1611		
1612	Pratyush Maini. 2023. Phi-1.5 model: A case of comparing apples to oranges?	1667
1613		1668
		1669
	Pratyush Maini, Skyler Seto, Richard Bai, David Grangier, Yizhe Zhang, and Navdeep Jaitly. 2024. Rephrasing the web: A recipe for compute and data-efficient language modeling . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14044–14072, Bangkok, Thailand. Association for Computational Linguistics.	1614
		1615
		1616
		1617
		1618
		1619
		1620
		1621
	Hieu Man and Thien Huu Nguyen. 2024. Counterfactual augmentation for robust authorship representation learning . In <i>Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR ’24</i> , page 2347–2351, New York, NY, USA. Association for Computing Machinery.	1622
		1623
		1624
		1625
		1626
		1627
		1628
	Nina Markl, Lauren Hall-Lew, and Catherine Lai. 2024. Language technologies as if people mattered: Centering communities in language technology development . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 10085–10099, Torino, Italia. ELRA and ICCL.	1629
		1630
		1631
		1632
		1633
		1634
		1635
	Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. 2018. UMAP: Uniform Manifold Approximation and Projection . <i>Journal of Open Source Software</i> , 3(29):861.	1636
		1637
		1638
		1639
	Miriam Meyerhoff. 2006. <i>Introducing Sociolinguistics</i> . Routledge, London, UK.	1640
		1641
	Alessio Miaschi, Dominique Brunato, Felice Dell’Orletta, and Giulia Venturi. 2020. Linguistic profiling of a neural language model . In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 745–756, Barcelona, Spain (Online). International Committee on Computational Linguistics.	1642
		1643
		1644
		1645
		1646
		1647
		1648
	Timothee Mickus and Maria Copot. 2024. Stranger than paradigms word embedding benchmarks don’t align with morphology . In <i>Proceedings of the Society for Computation in Linguistics 2024</i> , pages 173–189, Irvine, CA. Association for Computational Linguistics.	1649
		1650
		1651
		1652
		1653
		1654
	George K Mikros and Eleni K Argiri. 2007. Investigating topic influence in authorship attribution . In <i>SIGIR’07 International Workshop on Plagiarism Analysis, Authorship Identification, and Near-Duplicate Detection</i> .	1655
		1656
		1657
		1658
		1659
	Peter Millican. 2003. The Signature stylometric system . Web download. University of Oxford.	1660
		1661
	Moran Mizrahi, Guy Kaplan, Dan Malkin, Rotem Dror, Dafna Shahaf, and Gabriel Stanovsky. 2024. State of what art? A call for multi-prompt LLM evaluation . <i>Transactions of the Association for Computational Linguistics</i> , 12:933–949.	1662
		1663
		1664
		1665
		1666
	Frederick Mosteller and David L. Wallace. 1963. Inference in an authorship problem: A comparative study of discrimination methods applied to the authorship	1667
		1668
		1669

1670	of the disputed Federalist Papers. <i>Journal of the American Statistical Association</i> , 58(302):275–309.	1726
1671		1727
1672	Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2023. MTEB: Massive Text Embedding Benchmark . In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 2014–2037, Dubrovnik, Croatia. Association for Computational Linguistics.	1728
1673		1729
1674		1730
1675		1731
1676		1732
1677		1733
1678		1734
1679	Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candes, and Tatsunori Hashimoto. 2025. s1: Simple test-time scaling . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 20286–20332, Suzhou, China. Association for Computational Linguistics.	1735
1680		1736
1681		1737
1682		
1683		
1684		
1685		
1686		
1687	Ana Cristina Munaro, Renato Hübner Barcelos, Eliane Cristine Francisco Maffezzolli, João Pedro Santos Rodrigues, and Emerson Cabrera Paraiso. 2024. Does your style engage? Linguistic styles of influencers and digital consumer engagement on YouTube . <i>Computers in Human Behavior</i> , 156.	
1688		
1689		
1690		
1691		
1692		
1693	Tempestt Neal, Kalaivani Sundararajan, Aneez Fatima, Yiming Yan, Yingfei Xiang, and Damon Woodard. 2017. Surveying stylometry techniques and applications . <i>ACM Computing Surveys</i> , 50(6):86.	
1694		
1695		
1696		
1697	Dong Nguyen. 2025. Collaborative growth: When large language models meet sociolinguistics . <i>Language and Linguistics Compass</i> , 19(2):e70010.	
1698		
1699		
1700	Dong Nguyen, A. Seza Dođruöz, Carolyn P. Rosé, and Franciska de Jong. 2016. Computational sociolinguistics: A Survey . <i>Computational Linguistics</i> , 42(3):537–593.	
1701		
1702		
1703		
1704	Dong Nguyen, Rilana Gravel, Dolf Trieschnigg, and Theo Meder. 2013. “How old do you think I am?” A study of language and age in Twitter . In <i>Proceedings of the International AAAI Conference on Web and Social Media (ICWSM)</i> , pages 439–448, Cambridge, USA. Association for the Advancement of Artificial Intelligence.	
1705		
1706		
1707		
1708		
1709		
1710		
1711	Dong Nguyen and Jack Grieve. 2020. Do word embeddings capture spelling variation? In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 870–881, Barcelona, Spain (Online). International Committee on Computational Linguistics.	
1712		
1713		
1714		
1715		
1716		
1717	Dong Nguyen and Esther Ploeger. 2025. We need to measure data diversity in NLP — better and broader . In <i>Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing</i> , pages 8823–8832, Suzhou, China. Association for Computational Linguistics.	
1718		
1719		
1720		
1721		
1722		
1723	Dong Nguyen, Laura Rosseel, and Jack Grieve. 2021. On learning and representing social meaning in NLP: A sociolinguistic perspective . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 603–612, Online. Association for Computational Linguistics.	1726
1724		1727
1725		1728
		1729
		1730
		1731
		1732
		1733
		1734
		1735
		1736
		1737
		1738
		1739
		1740
		1741
		1742
		1743
		1744
		1745
		1746
		1747
		1748
		1749
		1750
		1751
		1752
		1753
		1754
		1755
		1756
		1757
		1758
		1759
		1760
		1761
		1762
		1763
		1764
		1765
		1766
		1767
		1768
		1769
		1770
		1771
		1772
		1773
		1774
		1775
		1776
		1777
		1778
		1779
		1780

1781	parallel examples . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 8662–8685, Albuquerque, New Mexico. Association for Computational Linguistics.	1837
1782		1838
1783		1839
1784		1840
1785		1841
1786		1842
1787	Fuchun Peng, Dale Schuurmans, Shaojun Wang, and Vlado Keselj. 2003. Language independent authorship attribution using character level language models . In <i>Proceedings of the Tenth Conference on European Chapter of the Association for Computational Linguistics - Volume 1</i> , EACL '03, page 267–274, USA. Association for Computational Linguistics.	1843
1788		1844
1789		1845
1790		1846
1791		1847
1792		1848
1793		1849
1794	James W. Pennebaker, Ryan L. Boyd, Kayla Jordan, and Kate Blackburn. 2015. <i>The Development and Psychometric Properties of LIWC2015</i> . University of Texas at Austin, Austin, USA.	1850
1795		1851
1796		1852
1797		1853
1798	Drexel University PSAL. 2013. JSAN—The integrated JStyle and Anonymouth package . The Privacy, Security and Automation Lab (PSAL) at Drexel University.	1854
1799		1855
1800		1856
1801		1857
1802	Fanchao Qi, Yangyi Chen, Xurui Zhang, Mukai Li, Zhiyuan Liu, and Maosong Sun. 2021. Mind the style of text! Adversarial and backdoor attacks based on text style transfer . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 4569–4580, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	1858
1803		1859
1804		1860
1805		1861
1806		1862
1807		1863
1808		1864
1809		1865
1810	Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza: A Python natural language processing toolkit for many human languages . <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> , pages 101–108.	1866
1811		1867
1812		1868
1813		1869
1814		1870
1815		1871
1816	Justin Qiu, Jiacheng Zhu, Ajay Patel, Marianna Apidianaki, and Chris Callison-Burch. 2025. mStyleDistance: Multilingual Style Embeddings and their Evaluation . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 16917–16931, Vienna, Austria. Association for Computational Linguistics.	1872
1817		1873
1818		1874
1819		1875
1820		1876
1821		1877
1822		1878
1823	Ella Rabinovich, Raj Nath Patel, Shachar Mirkin, Lucia Specia, and Shuly Wintner. 2017. Personalized machine translation: Preserving original author traits . In <i>Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers</i> , pages 1074–1084, Valencia, Spain. Association for Computational Linguistics.	1879
1824		1880
1825		1881
1826		1882
1827		1883
1828		1884
1829		1885
1830		1886
1831	Francisco Rangel, Paolo Rosso, Moshe Koppel, Efsthios Stamatatos, and Giacomo Inches. 2013. Overview of the author profiling task at PAN 2013 . In <i>Working Notes of Conference and Labs of the Evaluation Forum (CLEF)</i> , Valencia, Spain. CEUR Workshop Proceedings.	1887
1832		1888
1833		1889
1834		1890
1835		1891
1836		1892
		1893
		1894
		1895
	Sudha Rao and Joel Tetreault. 2018. Dear sir or madam, may I introduce the GY AFC dataset: Corpus, benchmarks and metrics for formality style transfer . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 129–140, New Orleans, Louisiana. Association for Computational Linguistics.	1896
		1897
		1898
		1899
		1900
		1901
		1902
		1903
		1904
		1905
		1906
		1907
		1908
		1909
		1910
		1911
		1912
		1913
		1914
		1915
		1916
		1917
		1918
		1919
		1920
		1921
		1922
		1923
		1924
		1925
		1926
		1927
		1928
		1929
		1930
		1931
		1932
		1933
		1934
		1935
		1936
		1937
		1938
		1939
		1940
		1941
		1942
		1943
		1944
		1945
		1946
		1947
		1948
		1949
		1950
		1951
		1952
		1953
		1954
		1955
		1956
		1957
		1958
		1959
		1960
		1961
		1962
		1963
		1964
		1965
		1966
		1967
		1968
		1969
		1970
		1971
		1972
		1973
		1974
		1975
		1976
		1977
		1978
		1979
		1980
		1981
		1982
		1983
		1984
		1985
		1986
		1987
		1988
		1989
		1990
		1991
		1992
		1993
		1994
		1995
		1996
		1997
		1998
		1999
		2000
		2001
		2002
		2003
		2004
		2005
		2006
		2007
		2008
		2009
		2010
		2011
		2012
		2013
		2014
		2015
		2016
		2017
		2018
		2019
		2020
		2021
		2022
		2023
		2024
		2025
		2026
		2027
		2028
		2029
		2030
		2031
		2032
		2033
		2034
		2035
		2036
		2037
		2038
		2039
		2040
		2041
		2042
		2043
		2044
		2045
		2046
		2047
		2048
		2049
		2050
		2051
		2052
		2053
		2054
		2055
		2056
		2057
		2058
		2059
		2060
		2061
		2062
		2063
		2064
		2065
		2066
		2067
		2068
		2069
		2070
		2071
		2072
		2073
		2074
		2075
		2076
		2077
		2078
		2079
		2080
		2081
		2082
		2083
		2084
		2085
		2086
		2087
		2088
		2089
		2090
		2091
		2092
		2093
		2094
		2095
		2096
		2097
		2098
		2099
		2100

2005	Jacob Tyo, Bhuwan Dhingra, and Zachary C. Lipton.	Minghao Wu and Alham Fikri Aji. 2025. Style over substance: Evaluation biases for large language models .	2060
2006	2022. On the state of the art in authorship attribution and authorship verification . <i>arXiv preprint</i> .	In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 297–312, Abu Dhabi, UAE. Association for Computational Linguistics.	2061
2007	ArXiv:2209.06869 [cs].		2062
2008			2063
2009	Adaku Uchendu, Thai Le, Kai Shu, and Dongwon Lee.	Linyi Yang, Yaoxian Song, Xuan Ren, Chenyang Lyu, Yidong Wang, Jingming Zhuo, Lingqiao Liu, Jindong Wang, Jennifer Foster, and Yue Zhang. 2023. Out-of-distribution generalization in natural language processing: Past, present, and future .	2064
2010	2020. Authorship attribution for neural text generation . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 8384–8395, Online. Association for Computational Linguistics.	In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 4533–4559, Singapore. Association for Computational Linguistics.	2065
2011			2066
2012			2067
2013			2068
2014			2069
2015	Suzanne Evans Wagner. 2025. Style and social meaning across the lifespan. <i>Connecting the Individual and the Community in Sociolinguistic Panel Research</i> , page 96.		2070
2016			2071
2017			2072
2018			2073
2019	Jan Philip Wahle, Terry Ruas, Yang Xu, and Bela Gipp.	Hanqing Zhang, Haolin Song, Shaoyu Li, Ming Zhou, and Dawei Song. 2023a. A survey of controllable text generation using transformer-based pre-trained language models . <i>ACM Computing Surveys</i> , 56(3):64:1–64:37.	2074
2020	2024. Paraphrase types elicit prompt engineering capabilities . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 11004–11033, Miami, Florida, USA. Association for Computational Linguistics.		2075
2021			2076
2022			2077
2023			2078
2024			2079
2025	Andrew Wang, Cristina Aggazzotti, Rebecca Kotula, Rafael Rivera Soto, Marcus Bishop, and Nicholas Andrews. 2023. Can authorship representation learning capture stylistic features? <i>Transactions of the Association for Computational Linguistics</i> , 11:1416–1431.	Jinghao Zhang, Yuting Liu, Wenjie Wang, Qiang Liu, Shu Wu, Liang Wang, and Tat-Seng Chua. 2025a. Personalized text generation with contrastive activation steering . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7128–7141, Vienna, Austria. Association for Computational Linguistics.	2080
2026			2081
2027			2082
2028			2083
2029			2084
2030	Janith Weerasinghe and Rachel Greenstadt. 2020. Feature vector difference based neural network and logistic regression models for authorship verification . In <i>Notebook for PAN at CLEF 2020</i> , volume 2695.	Yan Zhang, Zhaopeng Feng, Zhiyang Teng, Zuozhu Liu, and Haizhou Li. 2023b. How well do text embedding models understand syntax? In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9717–9728, Singapore. Association for Computational Linguistics.	2085
2031			2086
2032			2087
2033			2088
2034	Anna Wegmann and Dong Nguyen. 2021. Does it capture STEL? A modular, similarity-based linguistic style evaluation framework . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7109–7130, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.		2089
2035			2090
2036			2091
2037			2092
2038			2093
2039			2094
2040			2095
2041	Anna Wegmann, Dong Nguyen, and David Jurgens. 2025. Tokenization is sensitive to language variation . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 10958–10983, Vienna, Austria. Association for Computational Linguistics.	Zhehao Zhang, Ryan A. Rossi, Branislav Kveton, Yijia Shao, Diyi Yang, Hamed Zamani, Franck Dernoncourt, Joe Barrow, Tong Yu, Sungchul Kim, Ruiyi Zhang, Jiuxiang Gu, Tyler Derr, Hongjie Chen, Junda Wu, Xiang Chen, Zichao Wang, Subrata Mitra, Nedim Lipka, and 2 others. 2025b. Personalization of large language models: A survey . <i>Transactions on Machine Learning Research</i> .	2096
2042			2097
2043			2098
2044			2099
2045			2100
2046	Anna Wegmann, Marijn Schraagen, and Dong Nguyen. 2022. Same author or just same topic? Towards content-independent style representations . In <i>Proceedings of the 7th Workshop on Representation Learning for NLP</i> , pages 249–268, Dublin, Ireland. Association for Computational Linguistics.	Jiaxu Zhao, Meng Fang, Kun Zhang, and Mykola Pechenizkiy. 2025. Unmasking style sensitivity: A causal analysis of bias evaluation instability in large language models . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 16314–16338, Vienna, Austria. Association for Computational Linguistics.	2101
2047			2102
2048			2103
2049			2104
2050			2105
2051			2106
2052	E. Judith Weiner and William Labov. 1983. Constraints on the agentless passive . <i>Journal of Linguistics</i> , 19(1):29–58.	Hao Zheng and Mirella Lapata. 2022. Disentangled sequence to sequence learning for compositional generalization . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4256–4268, Dublin, Ireland. Association for Computational Linguistics.	2107
2053			2108
2054			2109
2055	Jennifer Williams and Simon King. 2019. Disentangling style factors from speaker representations . In <i>20th Annual Conference of the International Speech Communication Association: Crossroads of Speech and Language</i> , pages 3945–3949. ISCA.		2110
2056			2111
2057			2112
2058			2113
2059			2114
		Jian Zhu and David Jurgens. 2021. Idiosyncratic but not arbitrary: Learning idiolects in online registers	2115
			2116

reveals distinctive yet consistent individual styles. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 279–297, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Kangchen Zhu, Zhiliang Tian, Jingyu Wei, Ruifeng Luo, Yiping Song, and Xiaoguang Mao. 2024. *Style-Flow: Disentangle latent representations via normalizing flow for unsupervised text style transfer*. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 15384–15397, Torino, Italia. ELRA and ICCL.

Lal Zimman. 2019. *Trans self-identification and the language of neoliberal selfhood: Agency, power, and the limits of monologic discourse*. *International Journal of the Sociology of Language*, 2019(256):147–175.

Dimitrina Zlatkova, Daniel Kopev, Kristiyan Mitov, Atanas Atanasov, Momchil Hardalov, Ivan Koychev, and Preslav Nakov. 2018. *An ensemble-rich multi-aspect approach for robust style change detection*. In *Notebook for PAN at CLEF-2018*, page 3.

Chaoyuan Zuo, Yu Zhao, and Ritwik Banerjee. 2019. *Style change detection with feed-forward neural networks*. *Notebook for PAN at CLEF 2019*, 93.

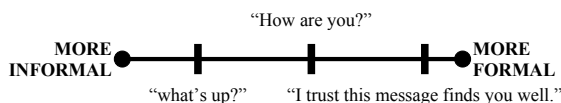


Figure 3: **Style is relative.** It might be more difficult or less interesting to make categorical judgments about a text’s style in isolation than, for example, judging if a text is more formal than another on a formality continuum. As Irvine (2001) writes on page 22, “It is seldom useful to examine a single style in isolation” and “attention must be directed to relationships among styles—to their contrasts, boundaries and commonalities.”

A Additional figures and tables

Fig. 4 provides a visual organization of the structure of this survey paper, Tab. 1 shows an overview of various predefined feature style operationalizations (§3.1), and Fig. 3 portrays an example of why style may require new solutions (§6).

B Motivating examples

B.1 Reasoning traces in the s1 dataset

We created Fig. 1 using the first 500 elements of the s1 datasets provided by Muennighoff et al. (2025) with reasoning traces generated by Gemini Flash Thinking Experimental and DeepSeek R1.¹² We

¹²“gemini_thinking_trajectory” and “deepseek_thinking_trajectory” column in <https://huggingface.co/datasets/simplescaling/s1k-1.1>

used a semantic representation model¹³ and a style representation model¹⁴ and UMAP (McInnes et al., 2018) with default settings.

Pioneering work found that the style of reasoning traces might be important to consider for the performance of reasoning models (Lippmann and Yang, 2025). Note, however, that their definition of style does not fully align with the definition used in this paper (e.g., including “non-linear thinking” as a style). In an ablation, we compare the semantic and style representations of the DeepSeek and Gemini teacher models and the distilled Qwen models on DeepSeek and Gemini. While the original Muennighoff et al. (2025) paper trains Qwen models only on Gemini reasoning traces, the authors later experimented with DeepSeek reasoning traces and found them to lead to better performance.¹⁵ We take the first 270 s1 reasoning traces as provided by Muennighoff et al. (2025) and use the fine-tuned Qwen models on Gemini¹⁶ and DeepSeek¹⁷ reasoning traces to generate reasoning traces¹⁸ for the first 270 Math500¹⁹ problems (Lightman et al., 2023). We use a different dataset from s1 to query student models to avoid artifacts of memorization. We choose Math500 as the distilled s1 Qwen models were also evaluated on it. See the results in Fig. 5 using UMAP visualization as before. We show that the style of the model distilled on Gemini reasoning traces is also closer in style to the Gemini reasoning traces than to the DeepSeek reasoning traces. Thus, the student model is effectively adopting the style of the teacher model (same for the DeepSeek model).

B.2 Rephrases of the MRPC dataset

Using synthetic data in pre- and post-training is increasingly common. We take the prompt from Maini et al. (2024) and use the Mistral-7B-Instruct-v0.1 model²⁰ (Jiang et al., 2023) to create Wikipedia-style rephrases of the first 500 elements of the MRPC dataset (Dolan and Brockett, 2005). We use the same models as in §B.1 for the semantic

¹³Hugging Face’s sentence-transformers/all-mpnet-base-v2

¹⁴Hugging Face’s AnnaWegmann/Style-Embedding

¹⁵<https://x.com/Muennighoff/status/1886405528777073134>

¹⁶<https://huggingface.co/simplescaling/s1-32B>

¹⁷<https://huggingface.co/simplescaling/s1.1-32B>

¹⁸By preceding the response with “\n<|im_start|>think\n”

¹⁹<https://huggingface.co/datasets/HuggingFaceH4/MATH-500>

²⁰<https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.1>

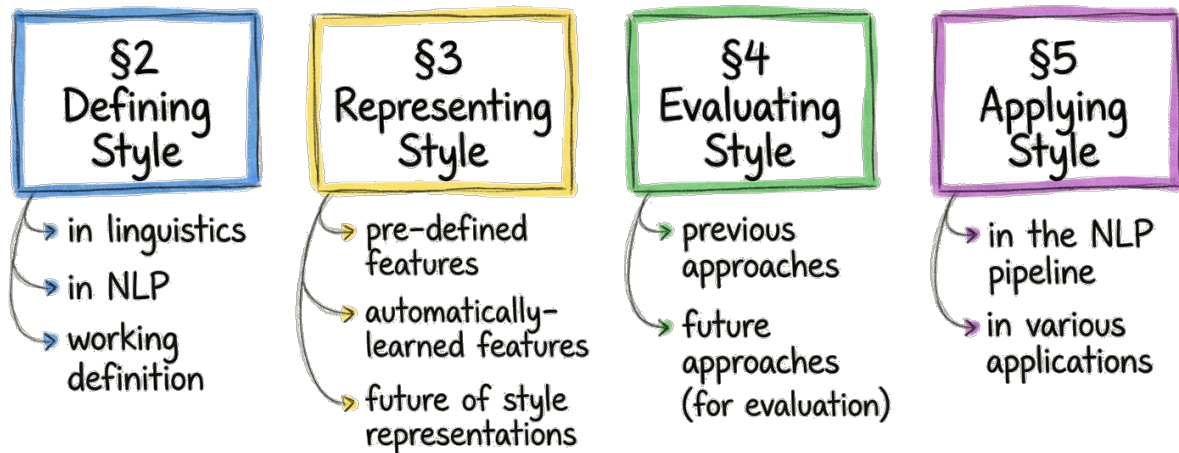


Figure 4: **Overview of the survey structure** This figure was digitalized from our own hand-drawn figure using NotebookLM and DALL-E. It keeps the same wording as the source material.

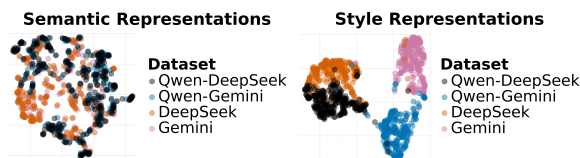


Figure 5: **Style representations of distilled Qwen models are close to teacher models** We compare reasoning traces on s1 for DeepSeek and Gemini models (Muenighoff et al., 2025) and reasoning traces on Math500 (Hendrycks et al., 2021) generated by models distilled on the s1 DeepSeek and Gemini reasoning traces respectively. The style representations (right) group the style of the student model closer to the style of the teacher model, while the semantic representations (left) overlap.

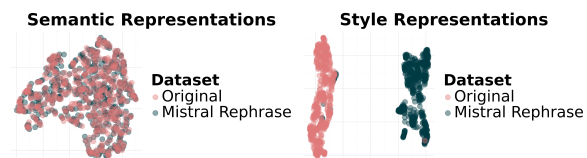


Figure 6: **Comparing semantic and style representations of LLM-rephrases** We compare MRPC sentences (Dolan and Brockett, 2005) and their LLM-generated “Wikipedia-style” rephrases using prompts from Maini et al. (2024). Style embedding models (right) can easily distinguish between the original and the LLM-rephrased sentences, while semantic embeddings (left) overlap. Studying stylistic diversity of LLM-rephrases is relevant as stylistic rephrasing is increasingly used in dataset curation for pre- and post-training.

and style representations as well as hyperparameters for the UMAP visualization. Style representations clearly distinguish the LLM rephrases from the original sentences, while semantic representations do not (Fig. 6).

B.3 Clustering writers of English by native language

We created Fig. 2 using the ETS Corpus of Non-Native Written English (LDC2014T06) (Blanchard et al., 2014). The corpus is comprised of English essays written by speakers of 11 non-English native languages as part of an international test of academic English proficiency, TOEFL (Test of English as a Foreign Language). We used LUAR²¹ (Rivera Soto et al., 2021), a style rep-

²¹<https://huggingface.co/rrivera1849/LUAR-MUD>

resentation model trained on the authorship verification task. Each point in the figure is an embedding of 5 TOEFL essays written by authors of the same native language picked at random. We reduce the dimensionality to two components using UMAP (McInnes et al., 2018) with default settings. Although the style representation was initially trained on the “idiolectal” authorship verification task (distinguishing authors based on their distinctive language use), Fig. 2 reveals that it also captures features indicative of the writer’s native language.

C Additions to style conceptualizations

Fully separating style and semantic meaning might be impossible.

Sociolinguists generally think of styles as different ways of saying the same thing. In every field that studies style seriously, however, this is not so.

— Penelope Eckert

A precise separation of semantic meaning and style poses practical challenges. It has been argued that, for example, only Labov (1972)’s original object of study—phonological variables—can leave semantic meaning untouched, whereas all other variables, including lexical and syntactic variables, will necessarily change the semantic meaning (Campbell-Kibler, 2011; Lavandera, 1978; Sun et al., 2023). Additionally, Eckert (2008, 2012) argues that using a certain style systematically connects an utterance to the social world, and that style thus influences social meaning. Others argue that any two forms must necessarily contrast in meaning (Clark, 1992). Some work in sociolinguistics sidesteps the problem of meaning equivalence by identifying and studying the contexts in which a set of linguistic forms are alternants without claiming equivalence (Campbell-Kibler, 2011; Christensen and Jensen, 2022). Nonetheless, we believe that attempting to separate style and semantic meaning has practical uses (see §2.2 or Weiner and Labov (1983)).

Style across research fields Several fields study linguistic style in some capacity. As discussed in the paper, *sociolinguistics* examines the relationship between language and society with a focus on language change and variation (Eckert, 2008; Labov, 1972). *Corpus linguistics* is the descriptive study of how language is actually used by analyzing text corpora (e.g. Biber, 1988; Biber and Conrad, 2019). Typical applications might include comparing language between different genres like scientific papers and news articles. *Forensic linguistics* involves the study of style in the context of law and crime investigation and is typically interested in recognizing a style or *idiolect* that helps distinguish an investigated individual (Coulthard, 2004). Practical insights from forensic linguistics also reciprocally influence *stylistics* and *stylometry*, which more generally study linguistic style in language. Stylometry applications include investigating the style of literary authors (Holmes, 1985) or attributing disputed literary works (Burrows, 2002; Mosteller and Wallace, 1963; Stamatatos, 2009). Style in *NLP* has been investigated to character-

ize authors (e.g., age or gender in Koppel et al., 2002; Nguyen et al., 2013), detect stylistic inconsistencies (Collins et al., 2004; Stamatatos, 2009), and adapt styles in machine translation (Niu et al., 2017, 2018; Rabinovich et al., 2017). Linguistic style also plays a significant role in related fields like *psycholinguistics*, or even in *communication* and *marketing*, such as by influencing consumer engagement (Munaro et al., 2024; ShabbirHusain et al., 2023) and purchases (Ludwig et al., 2013).

Note that these fields are not strictly separable. Methods from corpus linguistics can inform sociolinguistics, forensic linguistics can use methods from stylometry, and so on. Further, there are several fields that can be connected to linguistic style that we do not specifically discuss here, such as *discourse analysis*, *digital humanities*, *linguistic anthropology*, and *sociology*.

D Additions to representing style in NLP

Available predefined feature extraction tools

There are a multitude of tools available that automatically extract predefined features from text. The choice of tool and feature set, though, depends on various factors, such as preferred programming language, the nature of the data, and the goal of the task. Therefore, best practice is to systematically compare multiple feature sets, sometimes across tools, for each specific use case. Python tools include but are not limited to spaCy (Honnibal et al., 2020), Stanza (Qi et al., 2020), and NLTK (Bird et al., 2019) for general text processing, PAN submissions for authorship attribution (Weerasinghe and Greenstadt, 2020) and style change detection tasks (Strøm (2021), Zlatkova et al. (2018), LFTK (Lee and Lee, 2023) for extracting numerous stylometric features (but not n-grams), NeuroBiber and BiberPlus (Alkiek et al., 2025) for extracting Biber-style features, and StyloSpeaker (Aggazzotti and Smith, 2025) for extracting speech transcript features. Non-Python stylometric authorship tools include Stylo in R (Eder et al., 2016) and JStylo in Java (PSAL, 2013). Software that does not require programming includes LIWC (Pennebaker et al., 2015), which groups words into linguistically and psychologically meaningful categories; JGAAP (Juola et al., 2009) and Signature (Millican, 2003), which extract stylometric and n-gram features; and Coh-Metrix, which can measure more complex features like text cohesion (Graesser et al., 2004). We summarize these tools in Tab. 2.

Available automatically-learned models To the best of our knowledge, the available learned style representation models on HuggingFace are CISR²² (Wegmann et al., 2022), StyleDistance²³ (Patel et al., 2025), mStyleDistance²⁴ (Qiu et al., 2025), LUAR²⁵ (Rivera Soto et al., 2021) and Multilingual Style Representation²⁶ (Kim et al., 2025). Another model available via a private sharing site is LISA²⁷ (Patel et al., 2023). Following the discussion in §3.2, some style representations may capture more semantic features than others, and thus may prove to be more useful for different downstream tasks. We summarize these models in Tab. 3.

D.1 Additions to the future of style representations

? Automatic feature selection

Future work could attempt to create strategies to select predefined features that work best for different kinds of data and objects of study or develop an ensemble method that can select the best features dynamically.

? Including language modeling objectives

Previous work found that fine-tuning pretrained transformer models on style tasks can curiously lead to reduced performance on some style tasks compared to the pretrained base model (Patel et al., 2024; Wegmann and Nguyen, 2021). This might be connected to a difference in the object of study for the training and evaluation tasks. For example, using individuals as the object of study (e.g., using authorship verification as the training task) can lead to unlearning stylistic attributes that can vary for the same individual (e.g., the formality of their writing across online forums, job applications, and other contexts). When aiming to learn general-purpose style representations, it might be necessary to include further stylistic or continued language modeling objectives like masked language modeling.

²²<https://huggingface.co/AnnaWegmann/Style-Embedding>

²³<https://huggingface.co/StyleDistance/styledistance>

²⁴<https://huggingface.co/StyleDistance/mstyledistance>

²⁵<https://huggingface.co/rrivera1849/LUAR-MUD>

²⁶<https://huggingface.co/Blablalab/multilingual-style-representation-Llama-3.2>

²⁷<https://ajayp.app/posts/2023/11/learning-interpretable-embeddings-via-llms/>

? Improve content-independence

This was already mentioned in the main paper, but we highlight this point for more clarity again. “Generally, few style representations reach high scores on content-independence (🔧 App. Tab. 3) and might benefit from more exhaustive content disentanglement.”, see §4.1. Consider current content-disentanglement strategies in §3.2.

E Additions to evaluating style representations

Leverage measurement theory We give some concrete examples of how measurement theory (🔧 see Trochim et al., 2015) might be applied for style embeddings and benchmarks. Measurement theory can provide a theoretical framework that helps make sure different important validity and reliability aspects are considered in the evaluation of style representations and style benchmarks.

For style embeddings, *convergent validity* (i.e., does the measure show similar measurement for similar concepts?) might be assessed by testing that texts that have a similar style have similar representations. This could be done by perturbing texts in stylistically inconsequential ways (e.g., by swapping out named entities like “Maria has style.” to “Emma has style.”) and comparing their embeddings. *Discriminant validity* (i.e., Is the measure not sensitive to concepts it should not be related to?) might be assessed by confirming that texts that change in other aspects than style (e.g., content) are still embedded similarly. This has been assessed before by evaluating content-independence (§4). *Predictive validity* (i.e., Can the measure be used to predict something that it should be predictive of?) might be assessed by evaluating performance on downstream tasks that make use of style representations, such as style classification or style transfer tasks (§4). 🔧 See also Fang et al. (2022) for further inspiration.

For style benchmarks, *reliability* (i.e., Is the measure giving the same results with repeated measurement?) might be improved by making sure that the same seeds are used when applying the benchmarks—for example, when using style classification tasks and a classifier is trained on top of embeddings. 🔧 See also Bean et al. (2025) for further inspiration related to benchmark *validity*—for example, they suggest to employ sampling strategies like stratified sampling that are representative of the task space.

2414
2415

2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430

2431
2432
2433
2434
2435
2436
2437
2438
2439

2440
2441
2442
2443
2444
2445
2446

2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459

2460
2461
2462

F Additions to what style representations can enable

Disentangle internal representations It may be useful to disentangle LLM-internal representations of style to allow models to turn style information on or off as needed. Disentanglement approaches have helped cross-domain generalization (Yang et al., 2023; Zheng and Lapata, 2022) and might also help cross-style generalization. This can be especially relevant for stylistic tasks (e.g., machine text detection) that should rely on, and for semantic tasks (e.g., reasoning) that should not rely on, style information (Wegmann et al., 2025). Disentanglement might work especially well with mixture-of-experts approaches (Artetxe et al., 2022), with style-specific architectures (e.g., tokenizers) for relevant experts.

Authorship attribution Style representations can enable authorship verification and attribution tasks, including historical authorship attribution of disputed texts (Mosteller and Wallace, 1963), identifying harmful actors (Arabnezhad et al., 2020; Saxena et al., 2025), detecting plagiarism in educational contexts (Elkhatat et al., 2023), and attributing speakers from speech transcripts (Aggazzotti et al., 2024, 2025b; Tripto et al., 2023).

Considering style in annotations Human-written texts and labels can include spurious correlations as a result of annotation instructions (Gururangan et al., 2018). Style representations could be used to monitor the output of annotation efforts, and ultimately, to distinguish instructions that evoke more stylistically diverse annotations.

Bias identification and reduction As mentioned (§ 1), language models are often biased against certain styles, including those associated with marginalized groups. Approaches detailed in § 5.1, like curating training and evaluation datasets with more diverse styles, can improve performance across styles and thus reduce model bias. Further, it might be possible to use style representations to identify biased behavior of a trained model: For example, representations might be used to generate (§ 5.2) or cluster texts of similar styles, enabling systematic comparisons of model performances across style clusters.

Develop style measures With style measures we mean the broader class of methods and metrics that include style representations. One might, for exam-

ple, develop a metric that measures the formality of a text, returning values between 0 and 1. Style representations are similarly quantitative measures of stylistic properties, but they typically encode (latent) stylistic dimensions in a vector space. In this study, we focus on style representations, but they can be applied to develop style metrics.

F.1 Open questions in the application of style representations

We add open challenges in the application of style representations to different problems.

Circular evaluation in style transfer When performing generative tasks conditioned on style representations, such as authorship style transfer, difficulties can arise when comparing models. Various works (Horvitz et al., 2024a,b; Khan et al., 2024) train authorship style transfer models with the aid of style embedding models (§ 3.2) but also evaluate the adherence to the target style using style embedding models. When comparing two systems like ParaGuide (Horvitz et al., 2024a) and StyleMC (Khan et al., 2024), the former trained with CISR embeddings (Wegmann et al., 2022) and the latter with LUAR embeddings (Rivera Soto et al., 2021), it remains unclear which embedding space to use for evaluation without giving either model undue advantage. We encourage the community to investigate additional possibilities for evaluation (e.g., based on predefined features, cf. § 4.1) or establish a standard representation for training as well as evaluation.

Should we even care about styles for user-facing LLMs? Some recent work shows that more human-like outputs by LLMs might be dispreferred by humans and might lead to increased anthropomorphism (Cheng et al., 2025; Sandoval et al., 2025). This hints at a complex set of desiderata NLP researchers should consider when building LLMs and when using representations to steer LLMs toward generating texts in different styles. However, what style of output is preferred remains highly contextual (i.e., dependent on the setting) (Sandoval et al., 2025), and we believe that training on stylistically diverse corpora remains essential for LLMs to understand and engage with diverse human styles.

2463
2464
2465
2466
2467
2468
2469

2470
2471

2472
2473

2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493

2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508

2509
2510

2511
2512
2513

2514

2515
2516
2517
2518
2519

2520
2521
2522
2523
2524
2525
2526
2527
2528

2529
2530
2531
2532

2533
2534
2535
2536
2537
2538
2539
2540
2541
2542

2543

2544
2545
2546
2547
2548
2549

2550
2551
2552
2553

G Intended use and licenses for used artifacts

We only use models and datasets for motivating examples in our survey. We discuss their licenses and intended use below.

G.1 Datasets

s1k We use the s1k dataset provided by Muennighoff et al. (2025) and accessed at <https://huggingface.co/datasets/simplescaling/s1k-1.1>. The dataset was shared with an MIT license, which we adhere to.

MRPC We use the MRPC dataset provided by Dolan and Brockett (2005). The dataset is available on the Microsoft website at <https://www.microsoft.com/en-us/download/details.aspx?id=52398>. No license information is easily available. However, it is a widely used and shared dataset, and the paper mentions it is for the express purpose of stimulating research.

Math500 We use the Math500 dataset provided by Lightman et al. (2023). It was shared with an MIT license by OpenAI. See <https://github.com/openai/prm800k/>.

ETS Corpus of Non-Native Written English We use the ETS Corpus of Non-Native Written English (also known as TOEFL11 or LDC2014T06) provided by Blanchard et al. (2014). It is accessed via the Linguistic Data Consortium (LDC) at <https://catalog.ldc.upenn.edu/LDC2014T06>. The dataset is distributed under a specific LDC user license agreement restricted to non-commercial research use, which we adhere to.

G.2 Models

CISR We use Wegmann et al. (2022)’s CISR model at <https://huggingface.co/AnnaWegmann/Style-Embedding>. No clear license information is given, but the model was published in a research paper encouraging further use.

LUAR We use Rivera Soto et al. (2021)’s LUAR model at <https://huggingface.co/rrivera1849/LUAR-MUD>, shared with an Apache 2.0 license, which we adhere to.

SBERT We use an SBERT (Reimers and Gurevych, 2019) semantic representation model, <https://huggingface.co/sentence-transformers/all-mpnet-base-v2>, shared with an Apache 2.0 license, which we adhere to.

Mistral We use Jiang et al. (2023)’s Mistral-7B-Instruct-v0.1 model, <https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.1>. The model was shared with an Apache 2.0 license, which we adhere to.

s1 models We use Muennighoff et al. (2025)’s fine-tuned Qwen models on Gemini (<https://huggingface.co/simplescaling/s1-32B>) and DeepSeek (<https://huggingface.co/simplescaling/s1.1-32B>). Both models are shared with an Apache 2.0 license, which we adhere to.

H Identifying or offensive content in datasets

We use small existing datasets only for motivating examples (see §B). We do not release datasets. We do not expect the used datasets (§G.1) to include offensive content as they are reasoning datasets, crowd-worker created paraphrases and TOEFL essays. However, the TOEFL essays might include some personally identifying content. We did not take steps to remove identifiable cues or offensive content. We hope that the effect is negligible as the datasets were already publicly accessible and we only use them as motivating examples.

I Use of AI Assistants

We used ChatGPT, GitHub Copilot, and Claude Code for coding, to look up commands, and to generate individual functions for plotting. Generated functions were tested w.r.t. expected behavior. We used AI assistants (mostly Claude and ChatGPT) for concise rephrasing and grammatical error correction in writing. We used NotebookLM and DALL-E to generate one figure based on specific instructions including exact wording (see Appendix Fig. 4).

Type	Variable	Examples
PHONETIC	postvocalic /r/ intervocalic /t/ ...	more or less clear pronunciation of /r/ sound after vowel (Labov, 1972) full/flapped /t/ voicing between two vowel sounds (<i>writer</i> → <i>rider</i>) (Bell, 1984)
	MORPHO-LOGICAL	word endings nominalizations verb morphology ...
LEXICAL	word/token counts word/token ratios	number of words/tokens (Stamatatos, 2009) ratio of types to tokens, ratio of short/long words to token count, etc. (Altakrori et al., 2021)
	word/token n-grams word length sentence length vocabulary richness	for <i>n</i> of various lengths (Abbasi and Chen, 2008; Stamatatos, 2009) average word length (Biber, 1988), also cf. Grieve (2007) distribution of average sentence length, cf. Grieve (2007) hapax (dis)legomena, Yule’s I/K, number of unique tokens (Abbasi and Chen, 2008; Stamatatos, 2009)
	function words	grammar-functioning words, e.g., <i>the</i> , <i>be</i> , <i>to</i> (Abbasi and Chen, 2008; Mosteller and Wallace, 1963; Stamatatos, 2009)
	pronoun use	word frequency distributions of 1st, 2nd,... person pronouns (Biber, 1988; Pennebaker et al., 2015)
	hedge words	<i>at about</i> , <i>something like</i> as hedges in Biber MDA features; <i>maybe</i> , <i>perhaps</i> in tentative dimension in LIWC
	quantifiers	<i>each</i> , <i>all</i> as quantifier words or <i>everybody</i> , <i>anybody</i> as quantifier pronouns (Biber, 1988)
	...	
	SYNTACTIC	POS counts POS n-grams passive voice subordination features negation invariant <i>be</i> zero copula ...
DISCOURSE	contraction use discourse particle readability compression ...	<i>can’t</i> vs. <i>cannot</i> (contractions list ¹ , Biber (1988)) <i>well</i> , <i>now</i> (Biber, 1988) Flesch Reading Ease, Flesch Kincaid Grade Level, etc. (Python’s textstat ²) train a compression model on one text and use it to estimate how similar in style another text is, cf. Stamatatos (2009)
	ORTHO-GRAPHIC	character types character n-grams lengthening number substitutions misspellings acronyms/abbreviations ...

¹ https://en.wikipedia.org/wiki/Wikipedia:List_of_English_contractions

² <https://pypi.org/project/textstat/>

³ https://en.wikipedia.org/wiki/Wikipedia:Lists_of_common_misspellings/For_machines

⁴ https://en.wikipedia.org/wiki/SMS_language

Table 1: **Overview of predefined feature style operationalizations used in different fields.** Specific linguistic features that have been used to operationalize style and examples of each are categorized by linguistic level: phonetic (i.e., pronunciation and sound patterns), morphological (i.e., word structure and inflection), lexical (i.e., word choice), syntactic (i.e., sentence structure), discourse (i.e., larger structure), and orthographic (i.e., spelling and punctuation). Note that the categorizations might overlap, e.g., *g*-dropping might also be considered an orthographic or phonological variable, and character n-grams might encode different morphemes. These features have been investigated separately (Campbell-Kibler, 2009) and collectively (e.g., Abbasi and Chen, 2008; Biber, 1988; Neal et al., 2017; Stamatatos, 2009). This table was inspired by and partially filled with elements from other tables of stylometric features in these and other sources. For further references and examples consider also Grieve (2007) and Biber (1988).

Tool	Original Purpose	Language / Platform	Type	Link
spaCy (Honnibal et al., 2020)	General text processing	Python	library	github.com/explosion/spaCy
Stanza (Qi et al., 2020)	General text processing	Python	library	https://github.com/stanfordnlp/stanza
NLTK (Bird et al., 2019)	General text processing	Python	library	github.com/nltk/nltk
PAN 2020 AV (Weerasinghe and Greenstadt, 2020)	and AV	Python	Task subm.	github.com/janithnw/pan2020_authorship_verification
PAN 2021 SCD (Strøm, 2021)	SCD	Python	Task subm.	github.com/eivistr/pan21-style-change-detection-stacking-ensemble
PAN 2019 SCD (Zuo et al., 2019)	SCD	Python	Task subm.	github.com/chzuo/PAN_2019
PAN 2018 SCD (Zlatkova et al., 2018)	SCD	Python	Task subm.	github.com/machinelearning-su/style-change-detection
LFTK (Lee and Lee, 2023)	Stylometric feature extraction (no n-grams)	Python	library	github.com/brucelee/lftk
BiberPlus (Alkiek et al., 2025)	Biber-style feature extraction	Python	library	github.com/davidjurgens/biberplus
NeuroBiber (Alkiek et al., 2025)	Biber-style feature extraction	HF	Model	huggingface.co/Blablalab/neurobiber
MAT (Nini, 2019)	Biber-style feature extraction	Python	library	github.com/andreanini/multidimensionalanalysisstagger
StyloSpeaker (Aggazzotti and Smith, 2025)	Speech transcript feature extraction	Python	library	github.com/caggazzotti/styloSpeaker
Stylo (R) (Eder et al., 2016)	Stylometric authorship analysis	R	library	github.com/computationalstylistics/stylo
JStylo (Java) (PSAL, 2013)	Stylometric authorship analysis	Java	App	github.com/psal/jstylo
LIWC (Pennebaker et al., 2015)	Ling./psych. categories	SW (GUI)	App	www.liwc.app/
JGAAP (Juola et al., 2009)	Stylometric + n-gram features	SW (GUI)	App	evllabs.github.io/JGAAP/
Signature (Millican, 2003)	Stylometric + n-gram features	SW (GUI)	App	www.philocomp.net/texts/signature.htm
Coh-Metrix (Graesser et al., 2004)	Text cohesion and discourse features	SW (GUI)	App	soletlab.asu.edu/coh-metrix/

Table 2: **Comparison of common tools for extracting predefined features** The table summarizes their original purpose, programming language or platform, type of resource, and URL. Abbreviations: **AV** = authorship verification, **Task subm.** = shared-task submission, **SCD** = style change detection, **HF** = Hugging Face, **App** = standalone application, **SW** = non-programming software.

Model	Training Task	Languages	Content / Style Disentanglement	Interpretable?	Tasks
LUAR	AV	English	Weak	No	AR, MTD
CISR	AV	English	Medium	No	AV, MTD
StyleDistance	AV	English	Strong	No	AV, ST
mStyleDistance	AV	Multiple	Strong	No	AV, ST
LISA	AV	English	Strong	Yes	Unknown
Multilingual Style	AV	Multiple	Medium	No	AR, MTD

Table 3: **Comparison of open-source learned style representation models** The categorization is based on key dimensions including the languages supported, the measured strength of content/style disentanglement, interpretability, and the specific downstream tasks the models are have been found useful for. Note that the models may be useful for more tasks than stated here, the analysis is based on the authors’ experience with them. Acronym Definitions: **AR** = Authorship Retrieval, **AV** = Authorship Verification, **MTD** = Machine-Text Detection, **ST** = Style Transfer